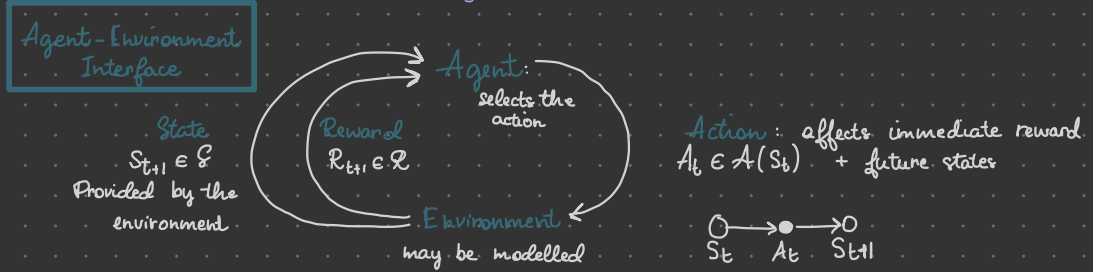


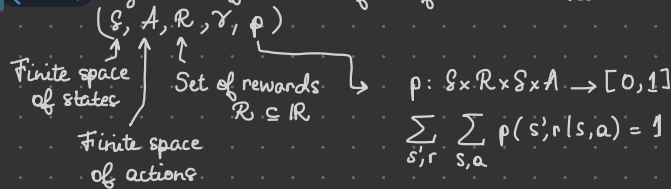
Reinforcement Learning Cheat Sheet

Fundamentals of RL



Goal: learn policy $\pi(a|s)$: map state \mapsto action taken, to maximise returns
 May be **deterministic**: $\pi(s) = a$; or **stochastic**: $\pi(a|s) \geq 0$, $\sum \pi(a|s) = 1$.

Markov Decision Process (MDP): general framework for sequential decision making, 5-tuple



Return: $G_t = \sum_{k=0}^H \gamma^k R_{t+k+1}$

Horizon may be ∞

Discount factor

$\begin{cases} \text{If } H \text{ finite: } \text{episodic task} \\ \quad \hookrightarrow \text{terminal state} \\ \text{If } H = +\infty: \text{continuing task} \\ \quad \hookrightarrow \text{no terminal state} \end{cases}$

$G_t = R_{t+1} + \gamma G_{t+1}$

Bellman's Equations

State-value function: $v_\pi(s) \equiv \mathbb{E}_\pi[G_t | S_t = s]$

Action-value function: $q_\pi(s, a) \equiv \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$

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

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) \times [r + \gamma v_\pi(s')]$$

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

Theorem: π_1 better than π_2 i.o.i. $v_{\pi_1}(s) \geq v_{\pi_2}(s) \forall s \in \mathcal{S}$.

There is always at least one optimal policy π_* (i.e. as good or better than any other).

$\pi_*(s) = \arg \max_a q_*(s, a)$

$$\Rightarrow q_*(s, a) = \max_{\pi} q_\pi(s, a), \quad v_*(s) = \max_a q_*(s, a)$$

Dynamic Programming (DP)

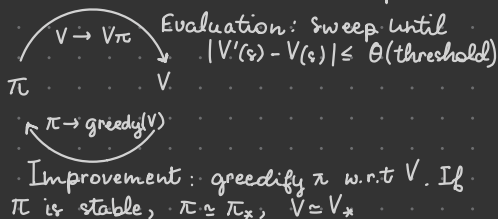
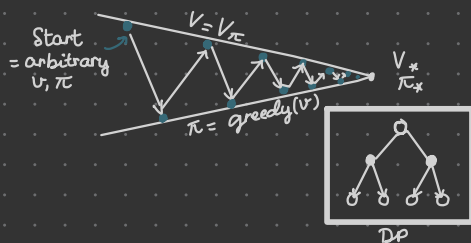
- **Bootstrapping**: use at least one estimated value in the update state for the same kind of expected value
- Requires a **model**.

Policy Evaluation (Prediction): Determine a policy's value function. **DP**: do both tasks iteratively.
Iteration (Control): Find a policy to maximize value functions.

Bellman's equations \mapsto update rule.
$$V'(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V(s)]$$

updated value old estimate

Policy Improvement Theorem: $V(s, \pi'(s)) \geq V(s, \pi(s)) \Rightarrow \pi' \geq \pi$. Greedified Policy = strict improvement unless already optimal.
 AND if $\exists s$ such that $q_\pi(s, \pi'(s)) > q_\pi(s, \pi(s)) \Rightarrow \pi' > \pi$.



General Policy Iteration (GPI): all algorithms iteratively performing evaluation + improvement.

Sample-based Learning Methods

Monte Carlo methods (MC)

- No bootstrapping
- Model-free

Relies on repeated random sampling. \mapsto averages sampled returns.



- Generate an episode according to π , recording $S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$
- From $T-1 \rightarrow 0$: $G_t = \gamma G_t + R_{t+1}$. Keep track of multiple observed returns.
 $V(S_t) = \text{average } [G_t] \text{ over generated episodes where } S_t \text{ was visited}$ or $Q(S_t, A_t) = \text{average } [G_t] \text{ over generated episodes where } S_t, A_t \text{ was visited}$

For non-stationary problems:
$$V'(S_t) = V(S_t) + \alpha [G_t - V(S_t)]$$

 α learning rate
 • To maintain exploration:
 ■ Exploring starts: A_0, S_0 picked at random.
 ■ ϵ -greedy: $\pi(a|S_t) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(S_t)|} & \text{if } a = \arg\max_a Q(S_t, a); \\ \frac{\epsilon}{|A(S_t)|} & \text{otherwise.} \end{cases}$

Off-Policy Learning

Use two policies: \rightarrow target $\pi(a|s)$ - evaluated policy, \rightarrow behavior $b(a|s)$ - used to generate episodes.

Introduce importance sampling ratio $\rho_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}$

$$V_\pi(s) = \mathbb{E}_b [\rho_{t:T-1} G_t | S_t = s]$$

Temporal Difference Learning (TD)

- Bootstrapping..
- Model-free.

$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] \quad \text{MC: uses samples to estimate } \mathbb{E}_{\pi}[\cdot]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \quad \text{TD: uses samples to estimate } \mathbb{E}_{\pi}[\cdot] \& v_{\pi}(S_{t+1})$$

TD update:

$$V'(S_t) = V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

TD-error δ_t



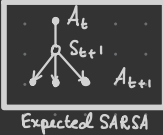
Sarsa:



$$Q'(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

On-policy
 A_{t+1} chosen from π before updating.

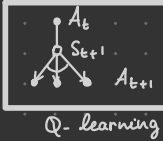
Expected Sarsa



$$Q'(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$$

Lower variance, more computationally expensive.
Can be used off-policy.

Q-learning



$$Q'(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

Special case of expected Sarsa

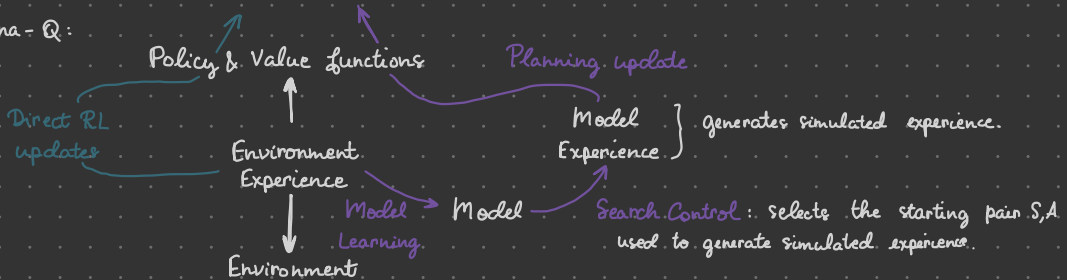
Planning

using a model to improve policy.

Q-planning:

1. Sampling: (S, A) picked at random, use sample model for S', R .
2. Q-learning update for (S, A) .
3. Greedify π w.r.t. $Q(S, A)$.

Dyna-Q:



To force exploration:

$$R \rightarrow R + \frac{k}{\sqrt{t}}$$

hyperparameter time since transition was last visited