

Reinforcement Learning

Reinforcement Learning - decision making

- ↳ No supervisor
- ↳ No instantaneous feedback
- ⇒ Agent gets action → influence environment
∴ influencing the data it can see

The RL Problem:

- $R_t \rightarrow$ Agent's goal is to maximise cumulative reward

Reward Hypothesis:

All goals \leftrightarrow maximisation of expected cumulative reward.

\Rightarrow Sequential decision making \rightarrow goal is one:

Time series of $A_t, o_t, R_t \rightarrow RL$

$$H_t = A_0, o_0, R_0, \dots, A_t, o_t, R_t$$

Now algorithm \rightarrow mapping from history to action

state \rightarrow info + determining what happens next

$$S_t = f(H_t)$$

$S_t^e \rightarrow$ info numbers to pick next obs/reward
 \hookrightarrow not seen by agent

$S_t^a \rightarrow$ info numbers to decide next action,
 \hookrightarrow used by RL algos.

$$S_t^a = f(H_t)$$

Information States

Markov property

$$P[S_{t+1} | s_t] = P[S_{t+1} | s_1, \dots, s_t]$$

$$H_{1:t} = (s_t) \rightarrow H_{1:t+1}$$

Future is independent of past, given the present
is enough to determine future.

$s_t^e \rightarrow$ Markov & H_t is markov (By defn)

env. state actually uses this state to

determine future actions = markov

What happens next depends on kind of state
that we will be in (not example)

Fully observable: $O_t = s_t^a = s_t^e$

Markov Decision Process (MDP)

Partially obs. state: agent state + env. state

Beliefs: $s_t^a = (P(s_t^e=s^1), P(s_t^e=s^n))$

(keeping prob distribution)

or RNN: $s_t^a = (s_t^a, \text{and } O_t(w))$

R.L agent could include following.

Policy: agent's behaviour fn:

- how agent picks its actions
- value fn: how good it is at particular state/action
- how much reward to expect.

Model: how agent thinks environment works

Policy: Deterministic policy $a = \pi(s)$

Stochastic policy $\pi(a|s) = P(A=a|s=s)$

Value: \rightarrow goodness / Badness of state

$$v_{\pi}(s) = E_{\pi}(R_t + \gamma R_{t+1} + \dots | s_t = s)$$



Behaviour that maximises this \rightarrow less enough risk.

Model: transition \rightarrow dynamics

rewards: - immediate reward

$$P_{ss'}^a = P(s' = s' | s=s, A=a)$$

$$R_s^a = E(R | s=s, A=a)$$
 on top prediction of reward

value Based:

value fn included

\Rightarrow No policy fn (Implicit understood)

how RL and planning problems are different

need to balance Exploration & Exploitation

find more info maybe by \rightarrow given info - maximise reward
losing reward

Lecture 2: i

Bandit - MDPs with 1 state!

$$P_{ss} = P \{ s_{t+1} = s' | s_t = s \}$$

↓
put 'i to n' to make a matrix

$$P = \begin{matrix} & \begin{matrix} P_{11} & \dots & P_{1n} \end{matrix} \\ \begin{matrix} \text{from} \\ \vdots \end{matrix} & \left[\begin{matrix} & & & \\ & \ddots & & \\ & & P_{nn} & \\ & & & \text{to} \end{matrix} \right] \\ \begin{matrix} P_{n1} \\ \dots \\ P_{nn} \end{matrix} & \begin{matrix} \text{from} \\ \vdots \\ \text{to} \end{matrix} \end{matrix}$$

Markov process (Markov chain) $\langle S, P \rangle$

$S = \{\text{set of all possible states}\}$

$P = \{\text{transition matrix}\}$

Sample - random sequence of these dynamics

Markov reward process

Reward fn: $R_s = E(R_{t+1} | s_t = s)$

says the immediate reward!

return $G_t \rightarrow \text{maximize goal / return}$

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

G_t is a random Sample \Rightarrow So no expectation

I chose to go to my office

I chose to go to night club

Discount factor

Uncertainty about rewards by taking
other choices!

Value function

$$v(s) = E(g_t | s_t = s)$$

What's next reward will you get until you
terminate?

Bellman equation

$$v(s) = E(R_{t+1} + \gamma v(s_{t+1}) | s_t = s)$$

$$= E(R_{t+1} + \gamma v(s_{t+1}) | s_t = s)$$

By law of expectation of expectation

$$v(s) = E(R_{t+1} + \gamma v(s_{t+1}) | s_t = s)$$

$$v(s) = R_t + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

$$V = R + \gamma P V \rightarrow \text{vector rep}$$

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

another level of complexity, i.e.

$$P_{ss'}^a = P[S_{t+1} = s' | s_t = s, A_t = a]$$

$$\pi(a|s) \subseteq P(A_t = a | S_t = s)$$

↑
This is what agent controls!!!

Policies.

MDP & MRP.

s_1, s_2, \dots, s_n is markov process (S, P^π)

s_1, r_2, s_2, \dots is markov reward process

(S, P^π, R^π, T)

where $P_{s,s'}^\pi = \sum_{a \in A} \pi(a|s) P_{s,s'}^a$

$$R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a$$

Action-value φ_π :

$$\varphi_{\pi_t}(s, a) = E_{\pi_t}(r_t + \gamma \varphi_{\pi_t}(s_{t+1}, A_{t+1}) |$$

$$\underbrace{\varphi_{\pi_t}(s, a) = E_{\pi_t}(r_{t+1} + \gamma \varphi_{\pi_t}(s_{t+1}, A_{t+1}) |}_{s_t = s, A_t = a}$$

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

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a V_{\pi}(s')$$

Q

$$V(s) = \sum_{a \in A} \pi(a|s) (R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a V_{\pi}(s'))$$

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a \sum_{a' \in A} \pi(a'|s') q_{\pi}(s', a')$$

$$V_{\pi} = R^{\pi} + \gamma P^{\pi} V_{\pi}$$

$$\pi_{\pi} = (1 - \gamma P^{\pi})^{-1} R^{\pi}$$

Optimal value function:

$$V_{*}(s) = \max_{\pi} V_{\pi}(s)$$

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



given action, over all policies

Solving MDP - finding q_{*}

Optimal Policy:

$$\pi \geq \pi' \quad \text{if} \quad V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s$$

Theorem: \exists an optimal policy π^*

$$\forall s \in \mathcal{S} \quad \forall a$$

$$V_{\pi^*}(s) = V_*(s)$$

$$q_{\pi^*}(s, a) = q_*(s, a)$$

deterministic optimal policy

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_{a \in A} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Bellman optimality eqn:

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

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_*(s')$$

$$V_*(s) = \max_a \left(R_s^a + \gamma \sum_{s' \in \mathcal{S}} \sum_{s \in \mathcal{S}} P_{ss'}^a V_*(s') \right)$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a \max_{a'} q_*(s', a')$$

Lecture 3:

Bellman eqn \rightarrow recursive decomposition

value fn \rightarrow stores / cache d. information

Planning in MDP:

for prediction

within

Input: MDP or MRP output v_t .

for control

Input: MDP output v_t & π_t .

Iterative Policy evaluation:

Synchronous backup:

~~update at iteration $t+1$ ($\forall s \in S$) ($v_{t+1}(s)$ from update)~~

For all iterations, for all 's', update $v_{t+1}(s)$ from $v_t(s')$

$s' \rightarrow$ successor state of s .

Bellman Expectation equation:

$$v^{t+1} = r^t + \gamma P^t v^t$$

Iterative policy:

$\pi_t \rightarrow$ Random policy

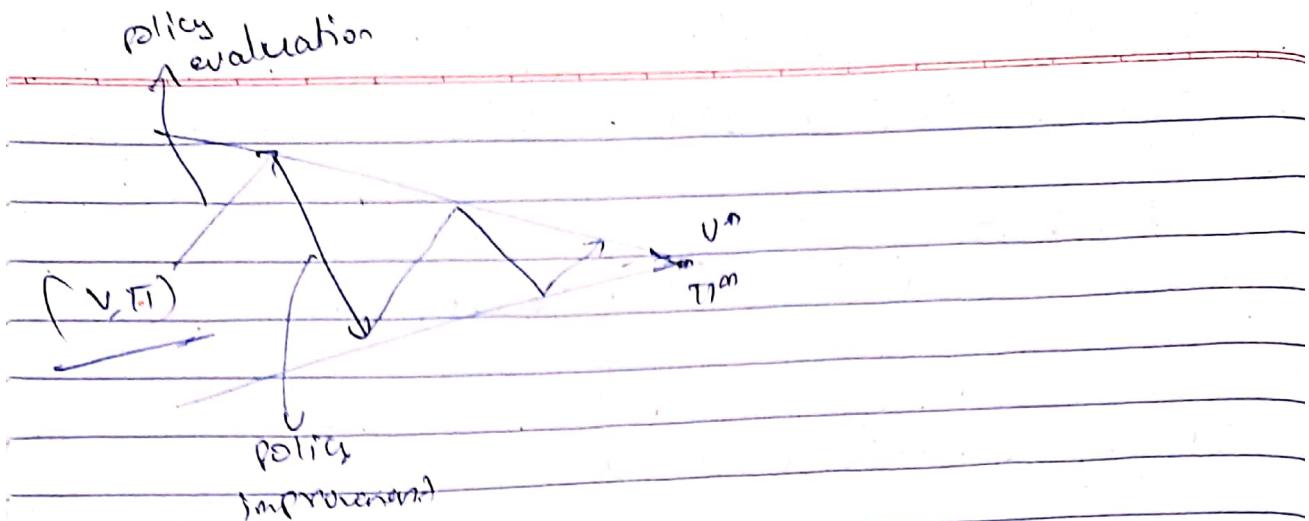
make actions not randomly
But using these policy.

a value fn \rightarrow figure out ~~optimal~~ optimal policy

Policing iteration:

Evaluate policy - $v_t(s) = \dots$

(π) Improve policy - choosing better action, i.e. greedy(v_t)



$$s \in S, \quad a = \pi(s)$$

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} q_{\pi}(s, a)$$

$$q_{\pi}(s, \pi'(s)) = \max_{a \in A} q_{\pi}(s, a) \geq q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

~~$v_{\pi'}(s) \Rightarrow v_{\pi}(s)$~~

$$v_{\pi'} \leq q_{\pi'}(s, \pi'(s))$$

$$= E_{\pi'} [R_{t+1} + \gamma v_{\pi'}(s_{t+1}) | s_t = s]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma q_{\pi'}(s_{t+1}, \pi'(s_{t+1})) | s_t = s]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots | s_t = s] \\ = v_{\pi'}(s)$$

If improvements stop.

$$q_{\pi'}(s, \pi'(s)) = \max_{a \in A} q_{\pi'}(s, a) = q_{\pi'}(s, \pi(s)) \\ = v_{\pi'}(s)$$

equalities of Bellman optimality

$$\Rightarrow \pi' = \pi$$

$\Rightarrow \pi$ is optimal policy.

Modified policy iteration:-

(\rightarrow improve only)

\rightarrow Consider only few iterations

Value iteration:-

Principle of optimality :-

\rightarrow Policy $\pi(a|s)$ achieves optimal value

$$V_{\pi}(s) = V_m(s) \text{ iff}$$

$\forall s'$ reachable from s

π^* achieves the optimal value from s .

$$V_{\pi^*}(s') = V_m(s')$$

$$V_m(s) \leftarrow \max_{a \in A} \left(R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a V_m(s') \right)$$

Use Bellmann optimality eqn as backup.

$$V_{\text{Bell}} = \max_{a \in D} (R_s^a + \gamma P_s^a V_K)$$

• state-value fns $V_{\pi}(s)$ or $V_m(s)$

\hookrightarrow complexity $O(mn^2)$ per iteration

m actions n states

• action-value function $q_{\pi^*}(s, a)$ or $q_m(s, a)$

\hookrightarrow complexity $O(m^2 n^2)$ per iteration

Asynchronous Dynamic programming :-

In place Dyn. prog.

Free Lecture 4 : Model-Predictive-control:-

Monte-Carlo learning:

Value = mean return!
Caveat: All episodes should terminate

$$V_{\pi} = E_{\pi} (G_t \mid s_t = j)$$

\uparrow terminates at a point!!
Used empirical mean Not at ' ∞ '!

Two types:-

(i) First visit (ii) Every visit

For each episode :-

$$s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_T$$

visit all / first state

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

$$S(s) \leftarrow S(s) + G_t$$

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

By central limit theorem, as $N(s) \rightarrow \infty$ $V(s) \rightarrow V_{\pi}(s)$

Incremental mean:

$$\mu_k = \mu_{k-1} + \frac{1}{N} (x_k - \mu_{k-1})$$

Defn of mean!

Incremental monte-carlo:

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

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

Sometimes we can use it?

Temporal difference learning:

guess \rightarrow guess \rightarrow ...

TD(0) algorithm:

update value $V(s_t)$ toward estimated return

$$\text{estimated return} = R_{t+1} + \gamma V(s_{t+1})$$

$$\therefore V(s_t) \leftarrow V(s_t) + \alpha \underbrace{(R_{t+1} + \gamma V(s_{t+1}) - V(s_t))}_{\Delta}$$

R_t

TD target

$\delta_t = \text{TD error}$

Bias / Variance Trade-Off:

Return $G_t = R_{t+1} + \dots + \gamma^{T-t} R_T \rightarrow$ unbiased estimate
of $V_T(s_t)$

$R_{t+1} + \gamma V_T(s_{t+1}) \rightarrow$ True target = unbiased estimate
of $V_T(s_t)$

$R_{t+1} + \gamma V(s_{t+1}) \rightarrow$ biased estimate of $V_T(s_t)$

\rightarrow can be used in random walk example.

Certainty Equivalence:

$$\sum_{k=1}^K \sum_{t=1}^{T_k} (g_t^k - V(s_t^k))^2$$

Solution to MDP (S, A, P, R, γ) that best fit data

$$P_{ss'} = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{T_k} I(s_t^k, a_t^k, s_{t+1}^k = s', a')$$

$$R_s^a = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{T_k} I(s_t^k, a_t^k = s, a) r_t^k$$

n step TD learning:

$$n=1 \quad (\text{TD}) \quad G_t^{(1)} = R_{t+1} + \gamma V(s_{t+1})$$

$$n=2 \quad \delta \quad G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+2})$$

$$n=\infty \quad (\text{MC}) \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots - \gamma^T R_T$$

Suppose

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(s_{t+n})$$

λ -step temporal difference:

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^{(n)} - V(s_t))$$

λ -rewards

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

weights $(1-\lambda)\lambda^{n-1}$

forward-view TD(λ)

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^\lambda - V(s_t))$$

Thm:-

Sum of offline updates is identical for forward-view & backward view TD(λ).

$$\text{let } \lambda \sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha (G_t^\lambda - V(s_t)) I(s_t=s)$$

$$TD(1) : E_x(s) = \gamma E_{x+1}(s) + r(s_t, s)$$

Total error by end of episode

$$= S_V = \gamma \delta_{V+1} + \gamma^2 \delta_{V+2} + \dots + \gamma^{T-1} \delta_T$$

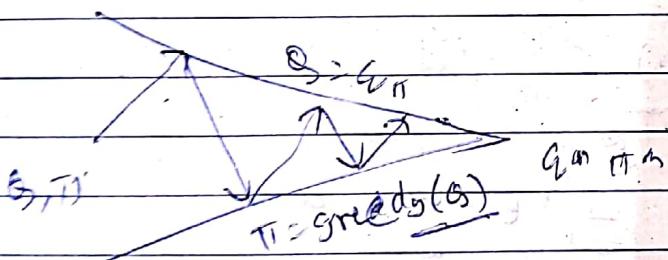
$\lambda = 1$, sum of TD errors

$$= G_t - V(s_t)$$

$$\text{for } \lambda \rightarrow \delta_t + (\gamma \lambda) \delta_{t+1} + (\gamma \lambda)^2 \delta_{t+2} + \dots$$

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \alpha \sum_{t=1}^T (\gamma \lambda)^{t-1} \delta_t = \alpha (G_V - V(s_V))$$

Reinforcement Model-Free Control



ϵ -Greedy exploration:

$$\pi(a|s) = \begin{cases} \frac{\epsilon}{m} + 1 - \epsilon & \text{if } a^* = \underset{a \in A}{\operatorname{argmax}} Q(s, a) \\ \frac{\epsilon}{m} & \text{otherwise} \end{cases}$$

Then:-

For any ϵ -greedy policy π , the ϵ -greedy policy π' with respect to Q_{π^*} is improved

$$V_{\pi'}(s) \geq V_{\pi}(s)$$

$$\begin{aligned}
 v_{\pi}(s, \pi'(s)) &= \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a) \\
 &= \frac{\epsilon}{m} \sum_{a \in A} a_{\pi'}(s, a) + (1-\epsilon) \max_{a \in A} q_{\pi}(s, a) \\
 &\geq V_{\pi}(s)
 \end{aligned}$$

Greediness in the limit with infinite exploration (GLIE)

$\pi' \Rightarrow \lim_{K \rightarrow \infty} N_K(s, a)$

$\dim \pi_K(a|s) = 1 \quad (a = \arg \max_{a' \in A} Q_K(s, a'))$

GLIE Monte Carlo ; $Q(s, a) \rightarrow q_m(s, a)$

For sarsa algorithm:

$Q(s, a) \rightarrow q_m(s, a) \quad \text{iff}$

$\pi_t(a|s) \quad \sum_{t=1}^{\infty} \alpha_t = \infty \quad \text{and}$

$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

n Step Sarsa:

$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (q_t^{(n)} - Q(s_t, a_t))$

$q_t^{(n)} = (1-\alpha) \sum_{n=1}^{\infty} \alpha^{n-1} q_t^{(n)}$

$\therefore Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (q_t^{(n)} - Q(s_t, a_t))$

Off-Policy learning:

$$\{s_1, a_1, r_2, \dots, s_T\} \sim \mu$$

$$\begin{aligned} E_{x \sim \pi}(f(x)) &= \sum p(x) f(x) \\ &= \sum \pi(a|x) \frac{p(x)}{\pi(a)} f(x) \\ &= E_{x \sim \pi} \left(\frac{p(x)}{\pi(a)} f(x) \right) \end{aligned}$$

$$q_t^{\pi_{old}} = \frac{\pi(A_t(s_t))}{\pi(A_{t+1}(s_{t+1}))} \frac{\pi(A_{t+1}(s_{t+1}))}{\pi(A_{t+2}(s_{t+2}))} \dots \frac{\pi(A_{T-1}(s_{T-1}))}{\pi(A_T(s_T))}$$

$$V(s_t) \leftarrow V(s_t) + \alpha (q_t^{\pi_{old}} - V(s_t))$$

Off-Policy control:

$$\pi(s_{t+1}) = \underset{a'}{\operatorname{argmax}} Q(s_{t+1}, a')$$

$$Q \text{ learning target} = R_{t+1} + \max_{a'} \gamma Q(s_{t+1}, a')$$