

REINFORCEMENT LEARNING

11

28 July.

- Sequential Decision Making - RL.

- Agent & Environment:

The agent interacts with an external environment. Agent has an internal state z_t , which it passes to its policy π to make a decision action $a_t = \pi(z_t)$.

Policy remains same, updating with each state.

Note: Tic Tac Toe - number of possible states, including X, O, & empty - including valid & invalid states - 19,683.

The environment responds with an observation o_{t+1} , which is used by the agent with the state update function $z_{t+1} = SUC(z_t, a_t, o_{t+1})$.

Reward: Immediate gain.

Goal: Maximize reward - get long-term reward

- Maximum expected utility principle:

$$V_\pi(s_0) = \mathbb{E}_p(a_0, s_1, a_1, \dots, a_T, s_T | s_0, \pi) \left[\sum_{t=0}^T R(s_t, a_t) \right] | s_0$$

where s_0 - agent's initial state.

$R(s_t, a_t)$ - reward function.

$V_\pi(s_0)$ - value function for policy π evaluated at s_0 .

Qn.

Draw a markov chain which represents a reinforcement learning.

29 July.

- State, Action, and Reward:
loops these 3.

Policy function π :

Decisions are taken based on the current state.

$$\pi(s, a) \mapsto [0, 1].$$

maps the pair (s, a) to a probability value.

$$\pi(a|s) = P(A=a | s=s).$$

\Rightarrow Probability of taking an action given a particular state.

Policy function can be a deterministic function or a random function.

Reward: After performing an action.

In Super Mario; the reward might be:

+1 - collecting coins.

+10,000 - winning/saving the princess.

-10,000 - Touching the turtles (?)

0 - Nothing is done.

Had it be +10000 for collecting coins & 1000 for saving princess, Mario would've been too busy going after coins.

State Transition: Old state $\xrightarrow{\text{action}}$ New state.

State transition can be random or deterministic.

Randomness to the state transition comes from the environment.

$$P(S' | S, a) = P(S'=s' | S=s, A=a).$$

P function outputs the probability of the new state

GOOD BOY



— / —

given the old state and the action.

- Two sources of Randomness:

One is from the action while the other is from the state.

→ Dot means all possible states.

$$A \sim \pi(\cdot | s).$$

Action is randomly sampled from the discrete probability distribution.

$$s' \sim p(\cdot | s, a).$$

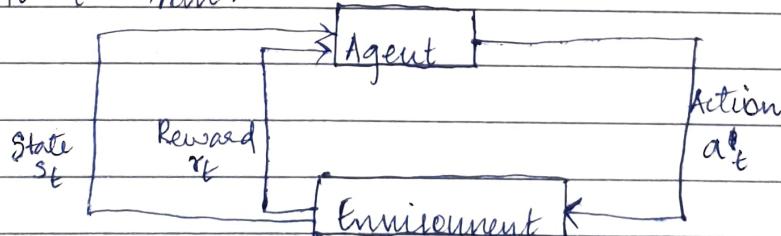
~ means 'sampled from'.

Randomness in Action is from the policy function.

Randomness in State Transition is from the state transition probability function.

Environment generates the new state by sampling the state transition function.

At t^{th} run:



- Play the game using AI:

- Observe the state s_t , select an action a_t

$$a_t \sim \pi(\cdot | s_t)$$

and execute a_t .

- AI will use the policy function to control the agent.

- The environment gives the new state s_{t+1} and



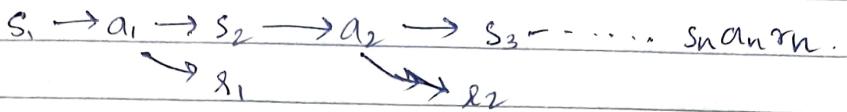
HEAD OVER HEELS

SPARKS FLY

SODA POPPS

— / / —

reward it :



$s_1 a_1 s_1, s_2 a_2 s_2, \dots, s_n a_n s_n$

episode

What is a good policy?

A good policy leads to a big cumulative reward.

$$\left[\sum_{t=1}^n \gamma^{t-1} r_t \right]$$

where γ (gamma) is the discounting factor.
 r_t is the reward.

31 July.

RL uses the rewards to guide learning of the policy function. Many algorithms are derived to learn good policy functions.

We record the trajectories as we train with samples.

We let the agent interact with the environment and record the trajectories.

- Rewards and Returns:

Return (cumulative future reward).

Let U_t be the return at time t .

$$U_t = R_t + R_{t+1} + \dots + R_n.$$

= sum of all future rewards from time t to the end.



— / /

Since future is uncertain, we introduce 'discounted return'.

Discounted return:

$$U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{n-t} R_n.$$

where γ (gamma) is the discounting factor, b/w the range $[0, 1]$.

The above is called the weighted sum of rewards from t to end.

At the end of the game, s_t, s_{t+1}, \dots , we can calculate the discounted return.

Randomness:

At any point of time ' t ', the reward R_t, R_{t+1}, \dots is random.

Reward R_i depends on State s_i and Action a_i .

State is random:

$$s_i \sim p(\cdot | s_{i-1}, a_{i-1}).$$

Action is also random:

$$a_i \sim \pi(\cdot | s_i)$$

where π is the policy.

Action is sampled from all the possibilities in the current state.

If either s_i or a_i is random, then R_i is also random. Because reward R_i depends on s_i & a_i .

U_t depends on R_t, R_{t+1}, \dots

In turn, U_t depends on all the states & actions at t .

When a game ends, we have all the returns and U_t is just a number.

- Action-Value Function:

$$Q_{\pi}(s, a)$$

At time t , U_t is random.

$$U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{n-t} R_n.$$

How is the current situation being evaluated?

- Expectation: $E[x] = \sum x p(x)$.

Take the expectation of U_t , by integrating out the randomness, we will get a real number that reflects how good the current situation is.

$$Q_{\pi}(s, a) = E[U_t | S_t = s_t, A_t = a_t]$$

It depends on the current state s_t and action a_t .

To compute the expectation, we need the probability density function (PDF) of S and A .

$Q_{\pi}(s, a)$ depends on the policy function.

$Q_{\pi}(s, a)$ is known as critic, since it evaluates the present performance of the agent.

$Q_{\pi}(s, a)$ depends on s_t, a_t, π and p .

$Q_{\pi}(s, a)$ is independent of s_{t+1}, a_{t+1} to s_n and a_{t+1} to a_n .

SIT DOWN

- BE HUMBLE -



-- / / --

- State - Value Function:

Assesses the how good a state is.

5 Aug.

I Don't Remember - It was an online class.

6 Aug.

- Model based Vs Model Free Approach:

Model based: We have exact knowledge of what we are trying to do., lack of uncertainty.
Eg: Tic Tac Toe.

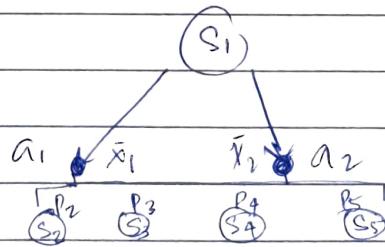
Model Free : Random moves; Not much info about
Eg: Mario.

Reinforcement learning.

Not much info about the env/agent.

Model based: Rewards from the env is known.
State transition probabilities are known.

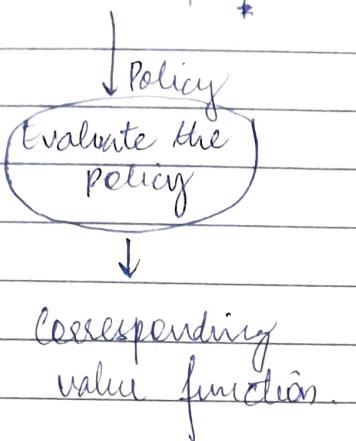
Eg:



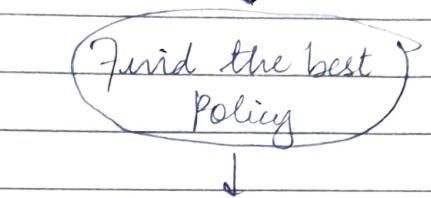
Model Free: Rewards from the env is unknown.
State transition probabilities are unknown.

Exploration has to be done, repeatedly
doing 'state-action-reward' loop.
Markovian approach?.

Predictive Problem



Control Problem (DYNAMIC)



- No input is provided.
- Goal: explore the policy space & find the optimal policy.

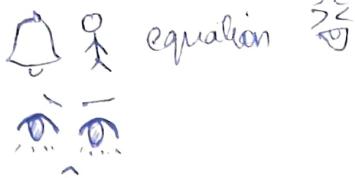
Classification of RL Algorithms:

	Predictive	Control
Model based planning	Dynamic programming- policy evaluation	- Dynamic prog - Value iteration- - DP policy iteration.
Model free RL	<ul style="list-style-type: none"> - Monte Carlo prediction. - Temporal Difference(TD). - TD(λ) backward. 	<ul style="list-style-type: none"> - Monte Carlo control. - SARSA . - SARSA backward - Q-learning . - Policy Gradient . - Actor Critic .

Interaction with the environment makes RL interesting!
???

Trajectory of interaction:

$s_1, a_1, r_1, s_2, a_2, \dots \rightarrow$ sequence can be in any order tho.
Eg: $s_1, a_1, s_3, a_3, s_7, a_7, \dots$ etc.



-- / /

- Bellman Equation:

Work backwards from a terminal state.



$R_8 \rightarrow$ Reward from taking an action to reach S_8 .

$S_1, a_1, r_1 \dots a_7, S_8, r_8$.

Return: All cumulative reward till that reward.



U_t - (cumulative) return - discounted cumulative sum of all previous rewards.

12 Aug

Did I miss any class? IDK 😐.

- Multi-Armed Bandits.

I can't focus at all, my head hurts because of someone's too-strong perfume. 😔

Regret = Difference b/w best possible expected reward & the actual reward we got.

$$\text{Regret}(T) = T p^* - \text{Reward}(T).$$

$p^* = \max_{i \in \{1, 2, \dots, k\}} p_i$, highest expected reward from one of the k arms.

$$\text{Avg regret}(T) = \frac{p^* - \text{Reward}(T)}{T}.$$

- Greedy (Naïve) Algorithm.
- Epsilon-Greedy Algorithm.
- Upper Confidence Bound (UCB) Problem.
- Thompson Sampling Algorithm.

20 Aug.

GAME THEORY

Missed about 15 minutes of the class.

I have very little idea of what's going on.

- Solutions:

Dominant Strategy - 2 types.

- Pure Nash equilibrium - prisoner's dilemma.
- ?
- Mixed strategy Nash equilibria.

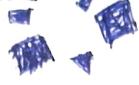
Matching pennies: Player 1 wins if both are heads or both are tails.

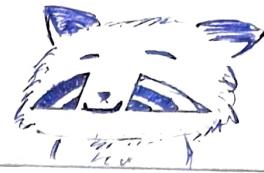
Player 2 wins if one head and one tail is observed.

- Games with no Nash Equilibria:

Finite set of players & finite strategy sets is needed for Nash equilibrium.

OOPS! Wrong Subject 公

- * check moodle for the lecture notes.
- * I missed at least 1 day per class.

 (car rental problem).



21 Aug.

- Planning by Dynamic Programming:

DP is used for optimization.

Problems \rightarrow Subproblems

For each subproblem, we try to find an optimal solution.

Done using recurrence.

- Markov Decision process.

A set of states an agent can explore.

Requirements for DP:

- Optimal Substructure: decomposing into subproblems.
 - Overlapping subproblems: subproblems reuse, cache solution to use again.
- * Markov decision satisfies both.

6 Oct

- Monte Carlo Method in RL:

Evaluation and Improvement (?) of policies (π):

- On-policy MC Method:

1. Generate experience using the current policy (π).
2. Update the policy π using experiences generated from above.
3. The policy is improved incrementally, and same policy is used to generate new experiences.

Eg:

Goal: Evaluate a target policy π , always choose action a in the given state to move to the next state.

Action $\in \{A, B\}$.

$S = \{S_1, S_2, \dots, S_n\}$. π

S	Action	Reward	Next State
s ₁	A	10	s ₂
s ₂	A	10	s ₁
s ₁	A	10	s ₃
:			

Estimate the value $v_{\pi}(s_1) = \sum_t r_t / N$,
where N - no. of experiences.

$$v_{\pi}(s_1) = \sum_i r_i / N = (10+10+\dots) / N.$$

Behaviour policy (β) .

- Off-policy MC method:
- 1. Generate exp or take actions to move from one state to another, using a different policy (β).
- 2. Estimate the value of target policy π using the experiences from above.
- 3. Target policy π is improved without seeing user to generate experiences.

Eg:

Same kind as previous example.

Actions $\in \{A, B\}$.

S	Action	Reward	Next State
s ₁	A	10	s ₂
s ₂	B	5	s ₁
s ₁	A	20	s ₃
:			

Importance Sampling : Identify The technique of identifying experiences relevant to π (target policy).

Weight = $P(A|S, \pi) / P(A|S, \beta)$.

[π :: the target policy says 'always choose A'].

$$\text{Var}(S_A) = 1/212 \approx 312 \approx 1.5.$$

$$\text{Weight} = 1/212 \approx 312 \approx 1.5.$$

$$\text{Var}(S_A) = (10 \times 1.5) + (5 \times 1.5) + (25 \times 1.5) + \dots / N.$$

9 Oct.

- Derivation of Importance Sampling:

$$\text{Weight} = \frac{P(A | S, \pi)}{P(A | S, \beta)} \cdot \frac{\pi}{\beta} = P.$$

$$x \leftarrow r.v.$$

Sampled from the probability distribution ' β ' and we want to estimate the expected value of x w.r.t the target distribution π .

$$E_{\pi}[x] = \sum_{x \in X} x \pi(x).$$

Important Sampling ratio ' P '.

$$\begin{aligned} E_{\pi}[x] &= \sum_{x \in X} x \pi(x) \\ &= \sum_{x \in X} x \pi(x) \frac{\beta(x)}{\beta(x)}. \end{aligned}$$

$$= \sum_{x \in X} x \underbrace{\pi(x) \beta(x)}_{\text{a new variable}}.$$

$$\rightarrow E_{\pi}[x]$$

$$\Rightarrow E_{\pi}[x] = \sum_{x \in X} x \pi(x) \beta(x).$$

PINOCCHIO

Assignment :- Gym? Jim? Use it.

- Try to implement RL MC or TD?

— / —

Now we can sample x from π and estimate its' expected value under π .

Here, we can change the policy if things don't work out, so getting more experience.

- Temporal Difference Learning (TD):

From the experiences, we "bootstrap", rather than waiting till the end of the episode.

Estimating how good the policy is on the go, like dynamic programming.

- Generate experiences and use it to improve the policy; but before the end of the episode.

Can be compared ^{with} to Monte Carlo - the opposite.

2 ways to do Bootstrap:

TD(0) update, TD(λ) update.

- TD(0) update:

It updates the value of the current state s_t based on the immediate reward R_{t+1} and the estimated value of next state $V(s_{t+1})$.

$$V(s_t) \leftarrow V(s_t) + \alpha [R_{t+1} + \underbrace{V(s_{t+1})}_{\text{TD-target}} - V(s_t)]$$



Value at a
given state



learning
rate

TD-target

~~?~~ Maximize \mathbb{E}
Tuoda venitkunni!!! ~~?~~

~~St~~

TD Error:

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

Goal is to minimize the error.

TD(0) \leftarrow TD(λ) \rightarrow MC.

[Current paying attention ~~?~~]

Optimization, control are the key points in RL.

- SARSA (On policy Control Algorithm):

It evaluates policy π that is used to take actions (behaviour). The update uses action A_{t+1} actually taken in the next state S_{t+1} .

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)].$$

E-Greedy policy can be used here (?).

$$S_{t+1} \leftarrow S_{B_t}$$

- Q-Learning (Off-Policy Control Algorithm):

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)].$$

13 Oct -

- Missed A Class. ~~??~~

LEARN DRIVING IN STRATEGY



14 Oct.

Function Approximation:

I'm too tired to write, sleepy.

- In RL, an agent explores the states & get rewards. This is stored in a look up table.
- Supervised learning.
- Gradient Descent ~~Algo~~ Method :

$$SGD: \quad w_{t+1} = w_t - \alpha \nabla L(w_t)$$

For TD update:

$$w_{t+1} = w_t - \alpha \underbrace{(\text{Target} - \hat{V}(s_t, w_t))}_{\text{TD error.}} \nabla \hat{V}(s_t, w_t).$$

Linear Function Approximation:

$$\text{Vector estimate } \hat{V}(s_t, w_t) = w^T X(s) = \sum_t w_t x_t(s).$$

$$\nabla \hat{V}(s_t, w_t) = x(s),$$

↳ feature vector.

TD(0) update:

$$w_{t+1} = w_t - \alpha [R_{t+1} + \gamma \hat{V}(s_{t+1}, w_t) - \hat{V}(s_t, w_t)] \cdot x(s_t)$$

Control Algorithm:

Find optimal policy π^* by approximating action value function Q .
(N times of prediction (state value)).



$$w_{t+1} = w_t + \alpha [R_{t+1} + \gamma \hat{q}(s_{t+1}, a_t) - \hat{q}(s_t, a_t, w_t)] \cdot \nabla \hat{q}(s_t, a_t, w_t)$$

• Q-Learning (Off Policy):

$$w_{t+1} = w_t + \alpha [R_{t+1} + \gamma \max_{a'} \hat{q}(s_{t+1}, a', w_t) - \hat{q}(s_t, a_t, w_t)] \cdot \nabla \hat{q}(s_t, a_t, w_t)$$