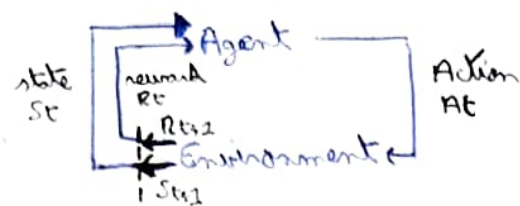


Focus on Reinforcement Learning



input: s_t at each time step

output: R_t at each time step

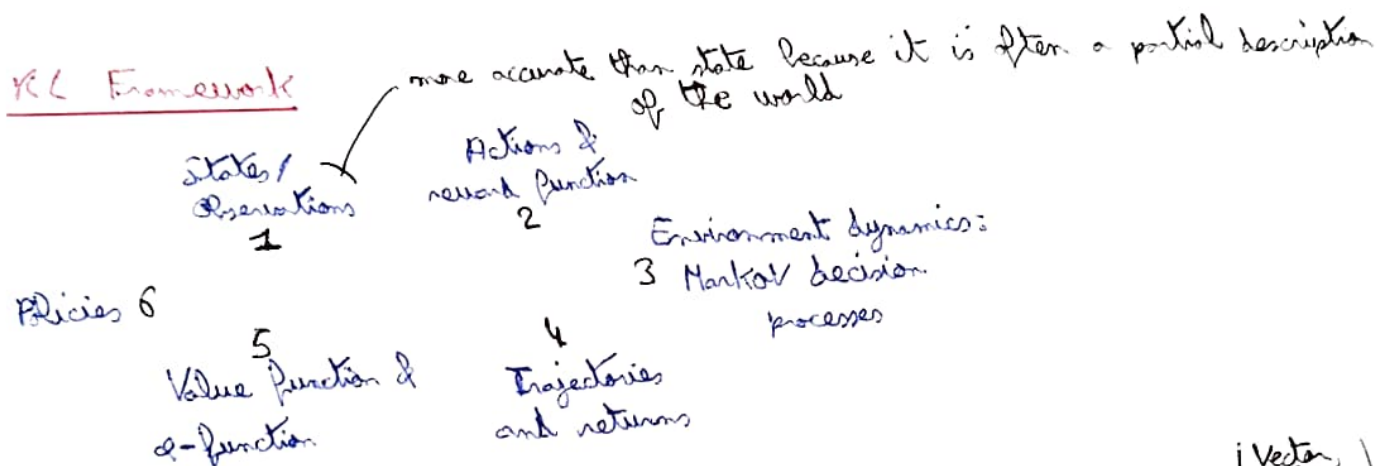
data: $(s_1, a_1, r_1, \dots, s_T, a_T, r_T)$

pick your own actions

goal: learn $\pi_\theta: s_t \rightarrow a_t$
to maximize J_π

Reinforcement Learning (RL): Learn some actions keeping on rewards and environment constraints

RL Framework



s_t : The state provide information for action decisions at timestep t (Vector, Matrix, Tensor, ...)

⌊ not the same info given depending on whether s_t is a vector or matrix or tensor (pong game example)

action a_t : at timestep t , the choices the agent can take
⌊ can be discrete or continuous (e.g. moving up or down pong paddle)

Let R be the reward function, then $R_t = R(s_t, a_t)$ a scalar indicating function (pong game ex.)
⌊ there is a reward for EACH action not just a group of actions

The environment dynamics of an RL model: Markov Decision Processes (MDP)

• Markov Decision Processes (MDP): discrete-time stochastic control processes with the Markov property

- Discrete-time: defines state the agent can be during problem-solving discrete
distinct points in time (no "in-between" states) → states are represented as

- Stochastic: Transition distribution: $p(s_{t+1} | s_t, a_t)$ is the proba to transition to s_{t+1} given the agent is in s_t and doing a_t
the agent may not end up in the intended state due to environmental randomness

- Markov Property: The state s_{t+1} depends only on the state s_t (even if the agent has been in any state)

⌊ a transition matrix for each action $A_i: p(s_{t+1} | s_t, a_t) = P_{s_t, a_t}$

want to apply RL? \rightarrow confirm that those 3 criteria have been met

- discrete action (up, down, or move)
- stochastic (the opponent introduces uncertainty in the point of view of the agent)
- Markov property (you move only from your current position to a new one)

A policy maps state to action \leftarrow Policy function as a NN for ex. (what to optimize are the weights of our NN)

\hookrightarrow deterministic policy: $a_t = \mu_\theta(s_t)$
 \hookrightarrow stochastic policy: $a_t \sim \pi_\theta(a_t | s_t)$

parameters θ are adjusted to optimize the agent's behavior

A trajectory refers to a sequence of states and actions: $\tau = (s_0, a_0, s_1, a_1, \dots)$

At the end of a trajectory, the agent accumulates rewards. This accumulated reward is called the return $R(\tau)$:

- Finite Horizon return ~~think~~ $R(\tau) = \sum_{t=0}^T r_t$ (limited number of timesteps)
 - Infinite Horizon return $R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$ (infinite or great number of timesteps)
- ($\gamma \in [0, 1]$ the discount factor \Rightarrow ($\gamma < 1$) future rewards \hookrightarrow valuable)
 \leftarrow ensures convergence

Value function at state s : $V_\pi(s) = E_{\tau \sim \pi} [R(\tau) | s_0 = s]$

$\hat{=}$ the reward we get of the trajectory τ that we expect to follow from policy π , given that $s_0 = s$

Q-function: $Q_\pi(s, a) = E_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a]$

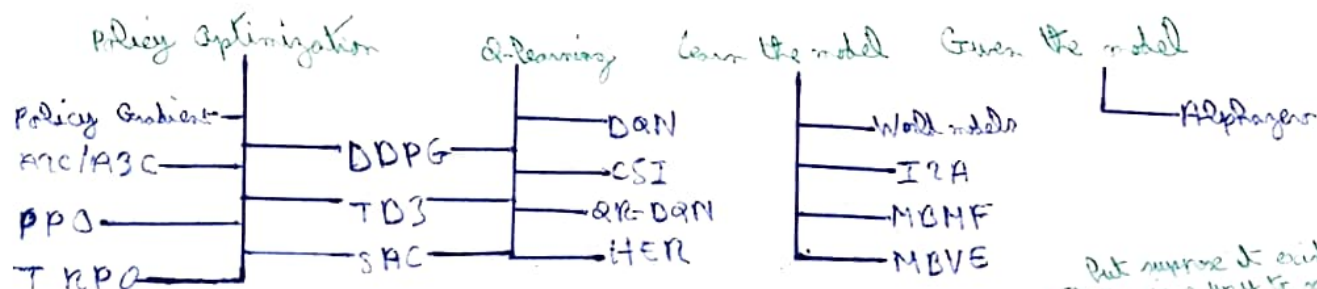
$\hat{=}$ the reward we get of the trajectory τ that we expect to follow from policy π , given that $s_0 = s$ and $a_0 = a$

a Taxonomy of RL Algorithms:

RL algos

Model-free RL

Model-based RL



Model-free RL: algos that does not use the transition probability $P(s_{t+1} | s_t, a_t)$ explicitly. They learn policies or value functions based on experience (trial and error)

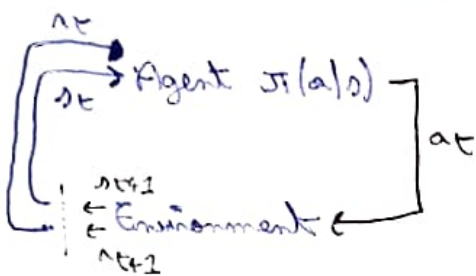
Model-based RL: algos that attempts to build a model of the environment, including $P(s_{t+1} | s_t, a_t)$, and uses this model to plan and optimize decisions

But suppose it exists and there is a way to compute it

Gradient Policy algorithms (PG);

- Optimize the policy directly (adjust its parameters to improve performance)
- train a policy to act based on observations

Goal: For the policy to choose actions that yield higher rewards



player in $s_t \rightarrow$ takes a_t based on $(\pi) \rightarrow$ the environment/opponent

(ex. of this type of agent: basketball player)

(stochastic case)
(ex.: coin flipping, learning with exploration)

adjust (π) after a trajectory (τ) using $R(\tau)$

reacts = player in s_{t+1}
 \downarrow
takes a_{t+1} based on s_{t+1}

\rightarrow often used for continuous action spaces where actions are sampled from a probability distribution rather than being discrete choices

$\hookrightarrow \pi(a|s) = \text{a distribution}$

state

Distribution Parameters

stochastic policy $\pi(a|s)$ as NN \Rightarrow state \rightarrow NN $\rightarrow (\mu, \sigma)$
 $\pi \sim N(\mu, \sigma)$

\hookrightarrow agent sample actions from the constructed distribution

\rightarrow then, gradient descent algorithm (ascent here in RL) for θ for $\pi(\theta)$, the performance measure of a policy $(\theta = \theta + \alpha \nabla_{\theta} J(\theta))$
 \uparrow LR = reward

RL Optimization Problem:

- The RL algo aims to maximize the expected return
- Goal: select the policy that achieves this maximization

Probability of a trajectory $\tau = (s_0, a_0, \dots, s_{T+1})$: $P(\tau|\pi_0) = P(s_0) \prod_{t=0}^T \pi_0(a_t|s_t) P(s_{t+1}|s_t, a_t)$

$$P(A \cap B) = P(A) \times P(B|A)$$

$$P(A \cap B \cap C) = P(A) \times P(B|A) \times P(C|A \cap B)$$

Expected return $J(\pi_0) = \int P(\tau|\pi_0) R(\tau) d\tau = \mathbb{E}_{\tau \sim \pi_0} [R(\tau)]$

Core RL problems: finding $\pi_0^* = \arg\max_{\pi_0} (J(\pi_0))$

Why $\nabla_{\theta} (J(\pi_0))$ and not $\nabla_{\pi_0} (J(\pi_0))$?
1- π_0 is a function: $J(\pi_0) = J(\pi(\theta))$
2- stochastic policy is itself parameterized by θ : $\pi = \pi(\theta) \Rightarrow$ (diff $\theta \Rightarrow$ diff π)

thanks to (2), $\nabla_{\theta} J(\pi_0) = \int \nabla_{\theta} P(\tau|\pi_0) R(\tau) d\tau$

$$\text{chain rule + log-trick + (1): } \nabla_{\theta} J(\pi_0) = \mathbb{E}_{\tau \sim \pi_0} [\nabla_{\theta} \log(P(\tau|\pi_0)) R(\tau)]$$

$$(1) + (2) \hookrightarrow \nabla_{\theta} J(\pi_0) = \mathbb{E}_{\tau \sim \pi_0} [\sum_{t=0}^T \nabla_{\theta} \log(\pi_0(a_t|s_t)) R(\tau)]$$

Estimation of $\nabla_{\theta} J(\pi_0)$: For $S = \{\tau^i\}_{i=1}^N$ with: N sampled trajectories

τ^i : trajectory of N trials from policy π_0

$$\hat{g} = \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log(\pi_0(a_t|s_t)) R(\tau^i) \right)$$

REINFORCE Algorithm :

1. sample $\{r^i\}$ from $\pi_\theta(a_t|s_t)$ (from the policy)

$$2. \nabla_\theta \bar{v}(\pi) = \sum_i \left(\sum_t \nabla_\theta \log(\pi_\theta(a_t|s_t)) \right) \underbrace{r(r^i)}_{\text{}} \left(\sum_t r(s_t^i, a_t^i) \right)$$

$$3. \theta \leftarrow \theta + \alpha \nabla_\theta \bar{v}(\pi)$$