

RL - unit 3

Reinforcement Learning Problem

It is about training an agent to make a sequence of decisions in an environment to maximize a cumulative reward. ~~Here's~~

Key Elements

- Agent.
- Environment.
- State.
- Action.
- Reward
- policy
- value function.
- Q -value

Agent

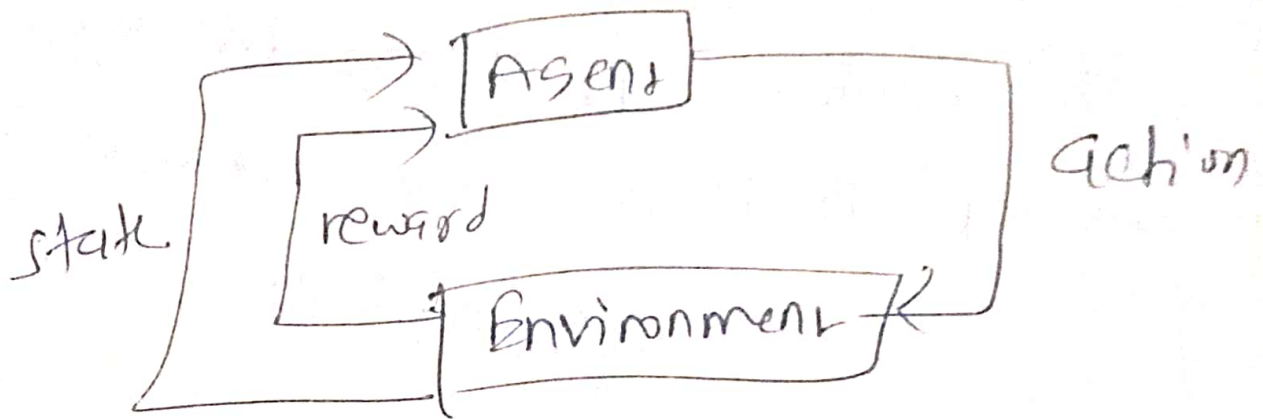
1. observe state
2. choose action

Action



Environment

3. provide Reward
4. New state

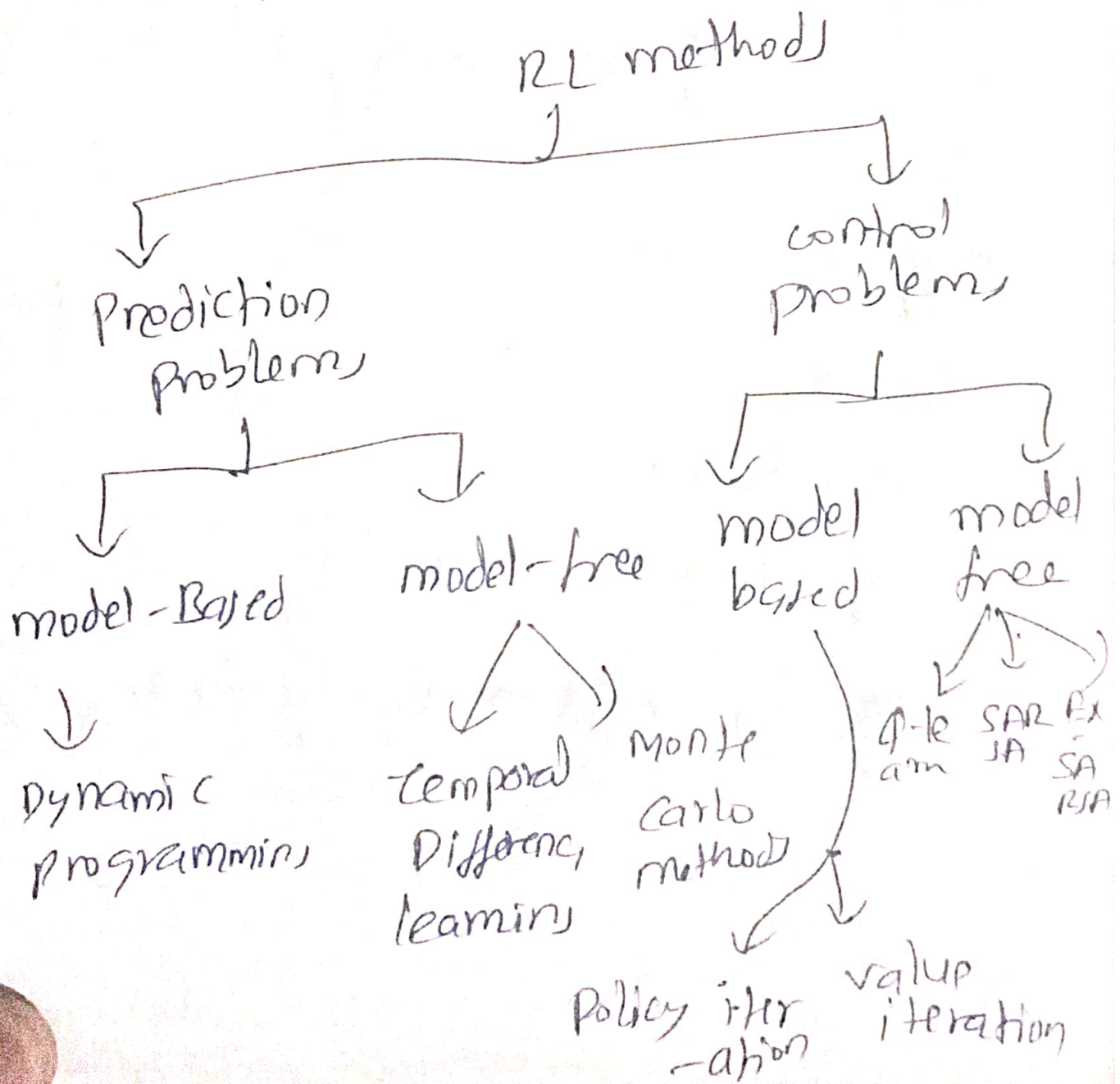


Application

- Game playing
- Robotics
- self driving car
- Recommendation system

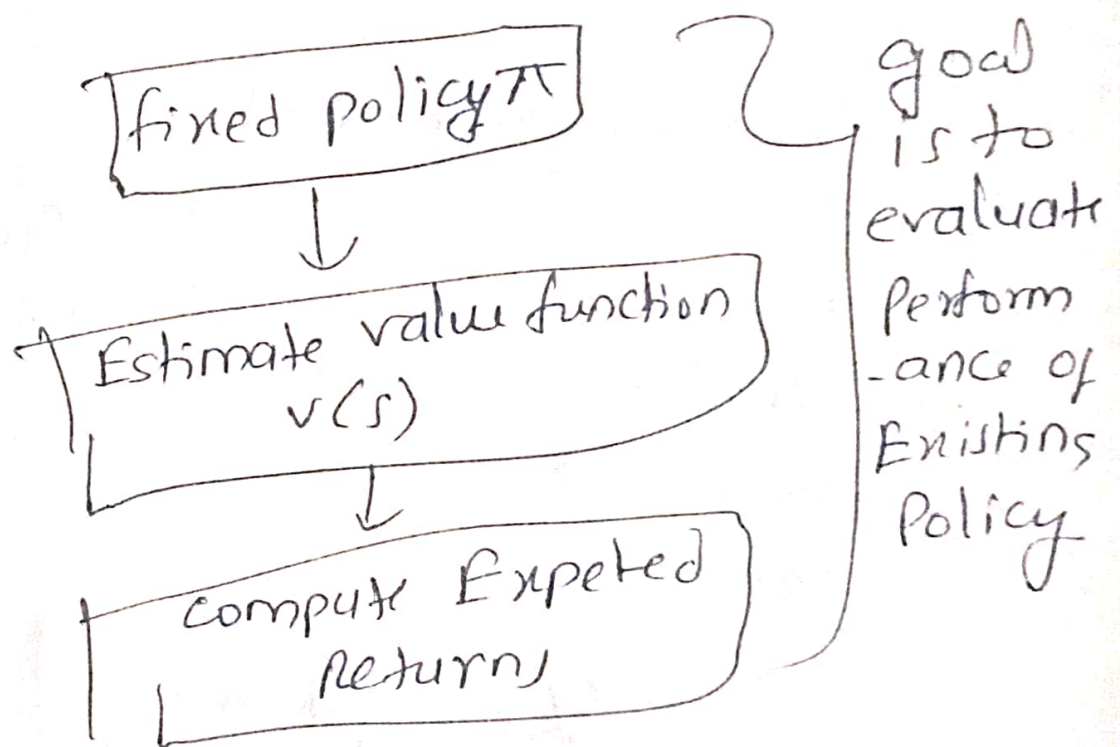
Prediction and control problems

RL is a subfield of ML where an agent learns to make decisions by interacting with an environment. The two fundamental problems in RL are prediction and control.



Prediction Problem

estimating the value function or the expected return for a given policy. The goal is to evaluate how good a particular policy is by predicting the future reward, the agent will receive.

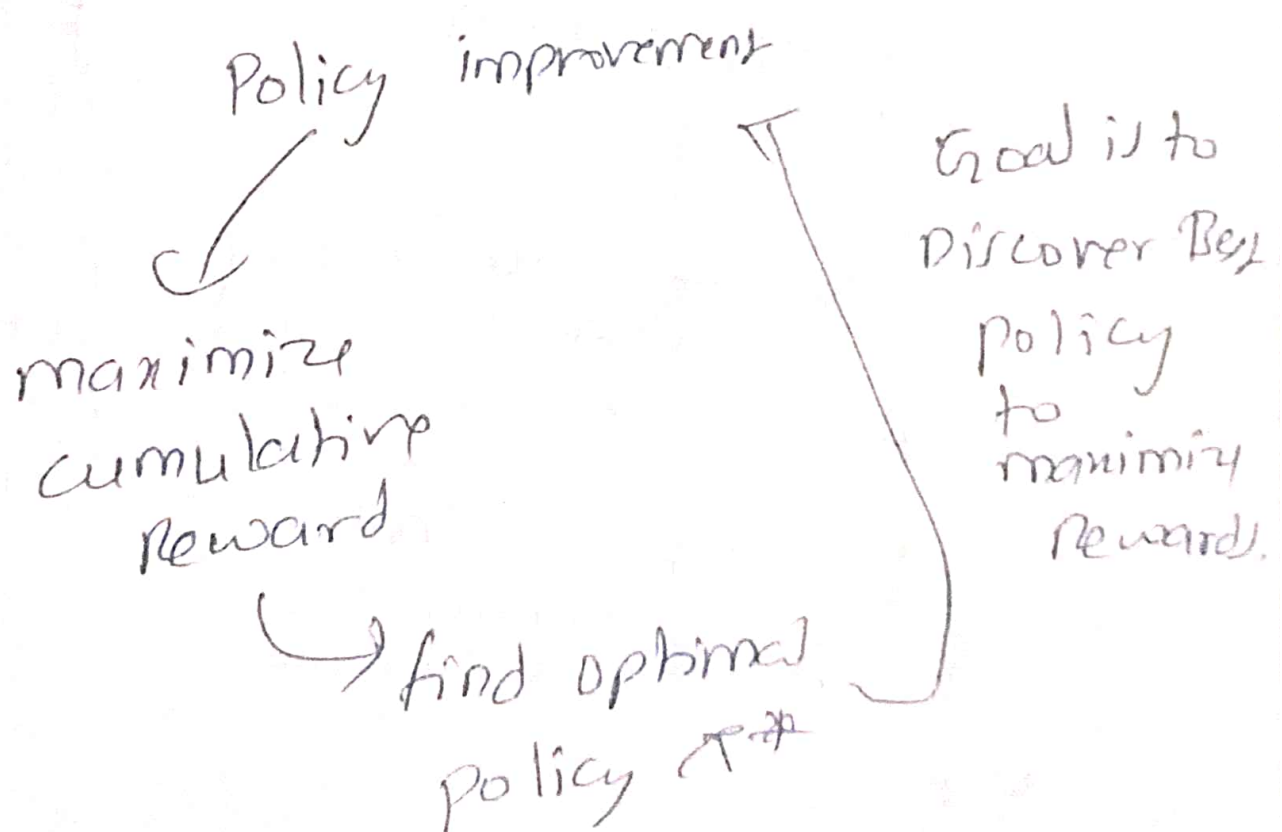


Ex:-

- weather forecasting
- Stock market
- Demand forecasting in economic.

Control problem

in RL involve finding the optimal policy π^* that maximizes the expected cumulative reward over time. The goal is to improve the policy iteratively to achieve the best possible performance.



Sys:-

- Autonomous vehicle navigation
- Industrial process optimization.

Prediction modeling

General
form:-

$$y = f(x) + \epsilon$$

y = predicted output

x = input features

ϵ → error term

Control system modeling

Feedback control
Equation:-

$$u(t) = k[r(t) - y(t)]$$

$u(t)$ → control input

$r(t)$ → Desired
reference state

$y(t)$ → current system
state

k → feedback gain.

Model based algorithms

These are class of methods that rely on learning a model of the environment to make decisions, unlike model-free algorithms, which directly learn a policy or value function without explicitly modeling the environment.

Key components

① → model learning

It learns a model of the environment,

→ Transition model

↳ Predict the next state given the current state and action

$$P(s' | s, a).$$

Reward model

↳ predicts the reward for a given state-action pair, $R(s, a)$.

→ The model is learned using supervised learning techniques, such as regression or neural networks.

② planning

once the model is learned, the agent uses it to simulate trajectories or plan actions.

→ planning methods include

↳ value Iteration or policy iteration

↳ model predictive control.

③ policy improvement

the agent improves its policy based on the planned actions.

learn Environment

Model

Simulate Trajectories

Optimal Policy

Planning

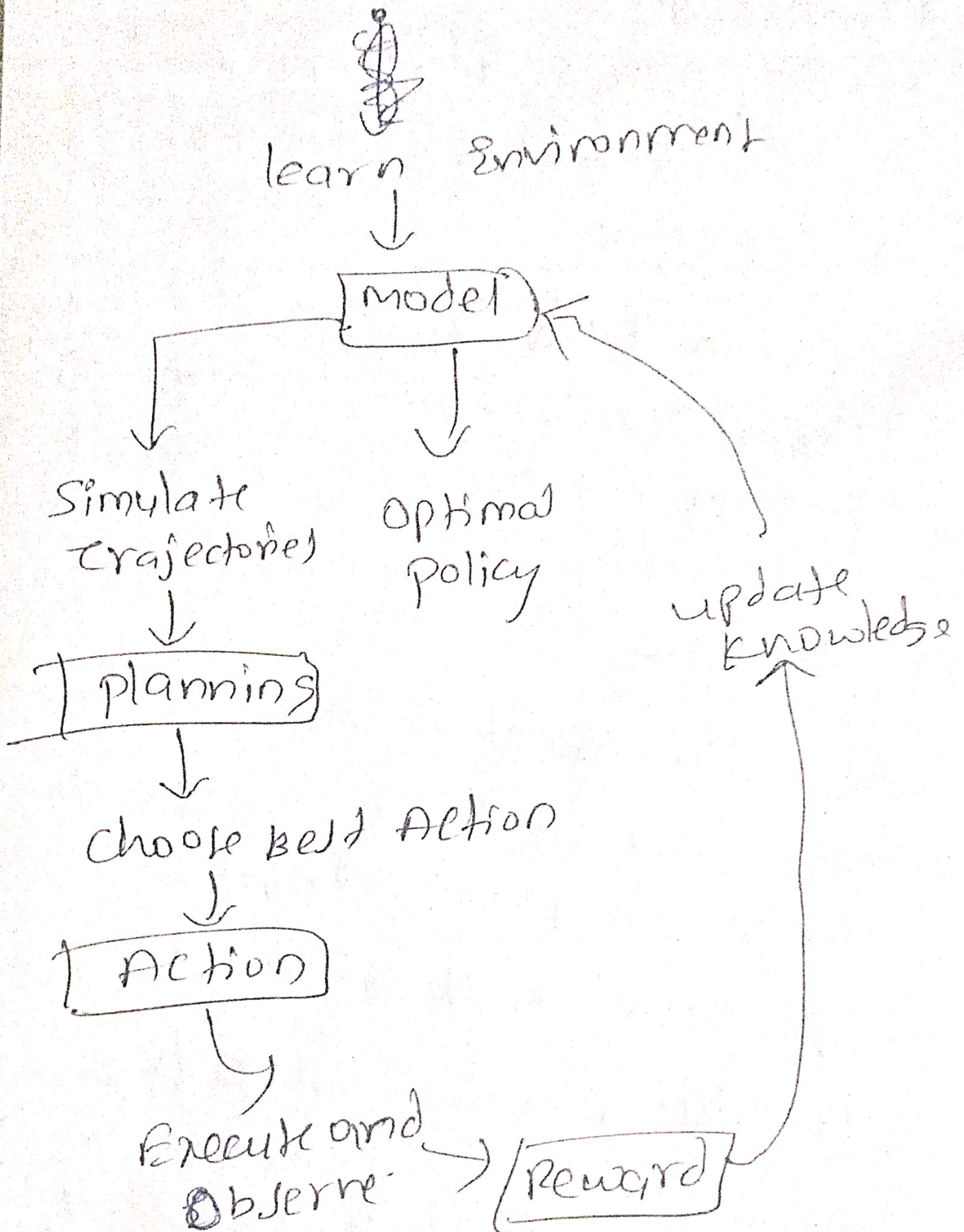
Choose Best Action

Action

Execute and Observe

Reward

update knowledge



Advantages

- sample Efficiency
- flexibility

Disadvantages

- computational cost
- complexity.

Ex:

- Dyna- ϕ .
- Model-Based policy optimization (MBPO).
- world models.

Monte carlo Methods for Prediction

→ this algorithms in RL is used for prediction and control.

They rely on sampling ~~complete~~ complete episodes of interaction with environment to estimate

value functions or policies.
→ Unlike TD methods, Monte Carlo require complete episode to compute returns, making them ~~inherent~~ inherently episodic.

Key concepts

① Episode.

Monte Carlo methods operate on complete episodes, which are sequence of state, action, rewards from initial to terminal state.

$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_T$

→ Episode should must terminate for monte carlo method work.

2. Return

The Return G_t is the total discounted reward from time step t onward.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

3. value function Estimation

The value function $v^\pi(s)$ is the expected return when starting in state s and following policy π .

$$v^\pi(s) = E[G_t | s_t = s]$$

Monte-Carlo prediction Alg

The goal is to estimate $v^\pi(s)$ for a given policy π .

① Initialize.

$V(s) \rightarrow$ A table to store the value function estimates for all state s .

+ Returns (r): ~~All s state of~~
↳ A list to store the returns
observed for each state s .

② Generate Episodes

↳ Use the policy π to interact
with the environment and generate
episode.

③ Compute Returns

for each episode, compute the
return G_t for each state s_t visited
in the episode.

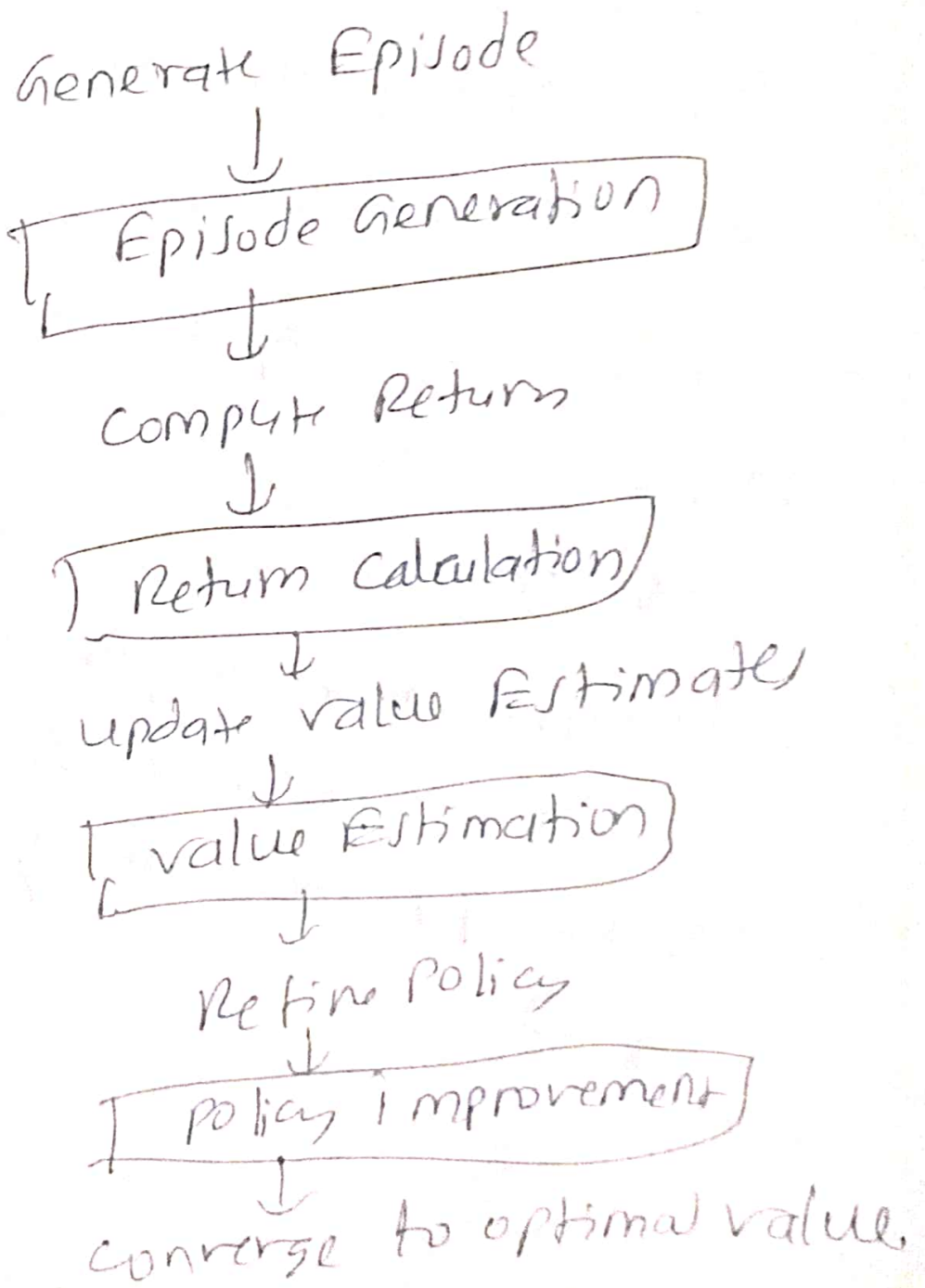
④ Update value functions

for each state s_t in the episode
append the return G_t to
 $\text{Return}(s_t)$

→ update the value estimate $v(s_t)$ as the average of all returns in $\text{Returns}(s_t)$.

⑤ Repeat

Repeat the process until the value estimate converges



online implementation of Monte carlo policy evaluation

Monte carlo policy evaluation is a method used to estimate the value function $v^\pi(s)$ for a given policy π by averaging the returns observed from sampled episodes.

→ The online implementation of monte . carlo policy evaluation updates the value function incrementally as new episodes are generated, rather than waiting through an episode in a batch

Key Step,

① Initialization

→ initialize the value function $V(s)$ for all state s arbitrarily
 $V(s) = 0$.

→ Initialize a counter $N(s)$ for each state s to keep track of the number of times the state has been visited.

② Generate Episode

↳ Simulate episode using the policy π , each episode is a sequence of states, actions, rewards.

③ Compute Return

for each state s in episode, compute the return G_t , which is the summation of discounted rewards from time t onwards.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

④ update value function

for each state s_t in the episode,
update the value function incrementally
using the observed return G_t !

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

⑤ Repeat

Repeat the process for multiple
episodes until the value function
 $v(s)$ converges to the true
value function $v^*(s)$.

First-visit
monte carlo

* updates each
state once per
episode

* unbiased,
higher variance

* lower (fewer
updates)

Even-visit
monte carlo

* updates each
state every time it
appears.

* slightly biased,
lower variance.

* Higher (more
updates).

Ex:

consider an episode with sequence

S_0, S_1, S_2, S_1, S_3

First-visit MC \rightarrow update $v(S_1)$ only for
the first occurrence of S_1 .

Even-visit MC \rightarrow update $v(S_1)$ for both
occurrences of S_1 .

Start Episode



[Episode]



Collect Returns



[Returns]



update $v(s)$



[value]



improve



[Policy]



converge

