

Chapter 1: INTRODUCTION: POLICY-BASED AND VALUE-BASED ALGORITHM

1/ Reinforcement Learning:

- Def: RL is concerned with solving sequential decision making problems.
- When solving question, we have objective and goals. We then take actions and get feedback from the world about how close we are to achieve the objective.
- RL problems can be expressed as a system consisting of an agent and environment.
 - + An environment produces information which describes the state of the system. This is known as a state.
 - + An agent interacts with the environment by observing the state and using information to select an action.
 - + The environment accepts the action and transition it into the next state.
 - + It then return the next state and a reward to the agent
 - + A cycle of state - action - reward complete, we said that 1 time step have passed/
 - + The cycle repeats until the env terminates.
- Policy = maps states to actions. An action will change an environment and affects what an agent observes and does next.
- RL problems have an object, **which is the sum of rewards received by an agent.**
- An agent's goal is to maximize the objectives by selecting the good actions.
- $(s,a,r) = (\text{state, action, rewards}) \Rightarrow$ This tuple at time t is called experience

2/ Reinforcement Learning as MDP:

- Transition function = MDP = math framework that models sequential decision making.
- To make environment transition function more practical, we turn it into a MDP by adding the assumption that transition to the next state only depends on the previous state and action
- The Markov property implies that the current state and action at time step t contain sufficient info to full
- The state:
 - + Is produced by an environment and observed by an agent.
 - + The state that is used b transition function.
- MDP formula for RL:
 - + S = set of state
 - + A = set of action
 - + $P(s_{t+1} | s_t, a_t)$ = state transition function of the environment
 - + $R(s_t, a_t, s_{t+1})$ = the reward function of the environment.

Algorithm 1 MDP control loop

```
1: Given an env (environment) and an agent:
2: for episode = 0, ..., MAX_EPISODE do
3:   state = env.reset()
4:   agent.reset()
5:   for t = 0, ..., T do
6:     action = agent.act(state)
7:     state, reward = env.step(action)
8:     agent.update(action, state, reward)
9:     if env.done() then
10:      break
11:   end if
12: end for
13: end for
```

3/ Learnable functions in RL:

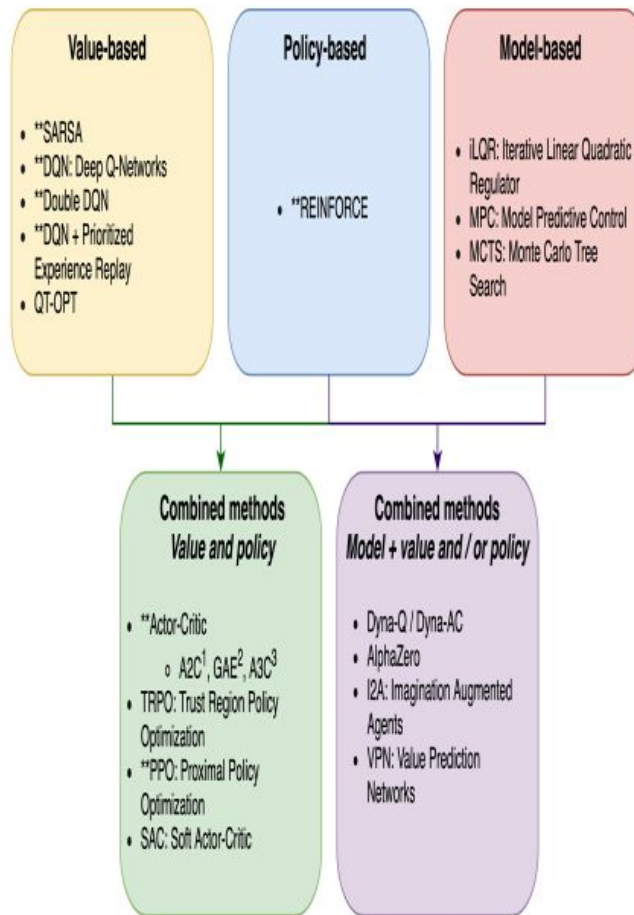
- Note: An agent can learn an action-producing function known as policy.
- There are properties of an environment that can be useful to an agent. 3 primary functions to learn in RL:

1. a policy, π , which maps state to action, $a \sim \pi(s)$
2. a value function, $V^\pi(s)$ or $Q^\pi(s, a)$, to estimate the expected return $E_\tau [R(\tau)]$
3. the environment model³, $P(s'|s, a)$

- + Policy = how an agent produce actions in the environment to maximise the objective
- + A policy can be written as $\pi(a|s)$ = probability of an action given a state s .
- + A value function provide info about the objective. It help an agent understand how good states and available action are in terms of expected future return.
- + A Q value function evaluate how good or bad a state-action pair is.

4/ Deep RL Algorithms:

- There are three major families of deep RL algorithms:
 - + Policy-based
 - + Value-based
 - + Model_based
- These model can be combine:



** : discussed in this book

1. A2C: Advantage Actor-Critic

2. A3C: Asynchronous Advantage Actor-Critic

3. GAE: Actor-Critic with Generalized Advantage Estimation

a/ Policy-based Algorithms:

- This algorithms learn a policy π . Good policies should generate action which produce trajectories that maximize an agent's object.
- If an agent needs to act in an environment, then it makes sense to learn a policy.
 - + What determine a good action at a given moment depends on the state
 - + Policy function take state s as input to produce an action a .
 - + Advantage: the algorithms are very general class of optimization methods.
 - + Disadvantage: very high variance and sample inefficient

b/ Value-based Algorithms:

- This algorithms learn V or Q function. It used the learned value function to evaluate (s, a) pairs and generate a policy.
- For ex: an agent's policy = select action a in state s with the highest estimated Q function.
- We use Q function better because it contain both state and action which is necessary to convert to policy.

- Value-based algorithms are typically more sample-efficient than policy-based algorithms.

c/ **Model-based Algorithms:**

- This algorithms either learn a model of an environment's transition dynamics or make use of a known dynamics model.
- Once an agent has a model of the environment, it can "imagine" what will happen in the future by predicting the trajectory of for a few time steps.
- An agent can complete many different trajectory predictions with different action sequences, then examine these different option to decide on the best action a to take.
- EX: Monte Carlo Tree Search.
- Model-based are very appealing since a perfect model can play out scenarios and understand the consequences of its action without having to actually act in the environment.
- However, the models are hard to come by. First, if the environment with large state space and action space can be very difficult to model. Second, the models are only useful when they can accurately predict the transition of the environment many step into the future.
- Model-based vs Model-free:
 - + Model-based = any algorithms that make use of the transition dynamics of an environment, whether learned or known in advance.
 - + Model-free algorithms = don't explicitly make use of the environment transition dynamics.

d/ **Combined Methods:**

- Actor-Critic algorithms = the policy acts and the value function critique the action.
 - + During training, a learned value function can provide a more informative feedback signal to a policy than the sequence of rewards available from the environment.
 - + During training, the policy learns using information provided by the learned value function

e/ **On-Policy and Off-Policy Algorithms:**

- Another distinction: An algo is on-policy **if it learned on the policy** ie, training can only utilize data generated from the current policy π_i .
 - + This implies that as training iterates through version of policy $\pi_1, \pi_2, \pi_3, \dots$ each training iteration only used the current policy that time to generate training data.
 - + All data must be discarded after training. This make on-policy methods sample inefficient, so they require more training data
 - + REINFORCE, SARSA, Actor-Critics
- Off-policy: Any collected data can be reused in training. This method is more sample efficient but may require a large memory to store data
 - + DQN, DDQN, PPO

f/ **SUMMARY:**

- Policy-based, value-based, model-based, or combined methods
- Model-based or model-free: whether an algorithms use a model of the environment's transition dynamics

- On-policy or off-policy: whether an algorithm learn with data gathered using just the current policy

5/ Deep Learning for Reinforcement Learning.