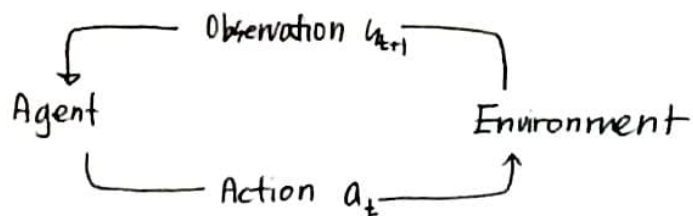


• Deep Reinforcement Learning

→ Basic Model



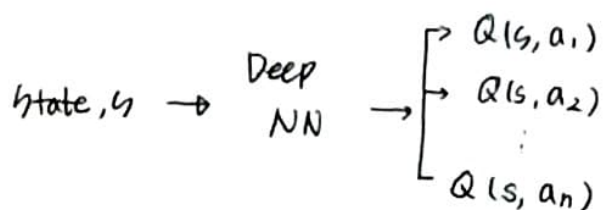
State $s_t \rightarrow$ action $a_t \rightarrow$ State $s_{t+1} \rightarrow$
 ↳ Reward r_t

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{i=t}^{\infty} \gamma^i r_i$$

$Q(s, a) = \mathbb{E}[R_t]$: expected total future reward in state s with certain a

$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$: policy that maximizes future reward

→ DQN

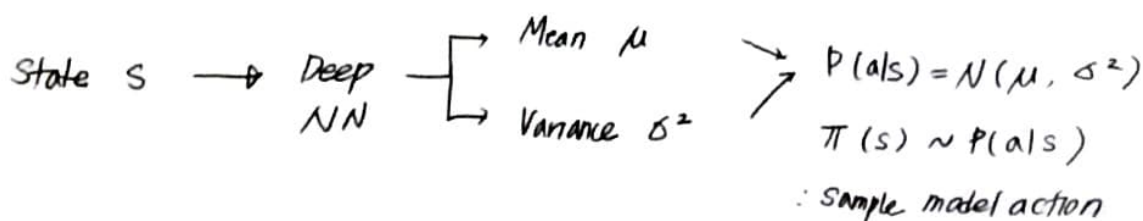


Maximize target return → train the agent

$$\mathcal{L} = \mathbb{E} \left[\left\| (r + \gamma \max_{a'} Q(s', a')) - Q(s, a) \right\|^2 \right]$$

but, cannot handle continuous action spaces

→ Policy Gradient: directly optimize the policy. Enabling modeling of continuous action space.



$$\mathcal{L} = \underbrace{-\log P(a_t | s_t)}_{\text{log likelihood}} \underbrace{R_t}_{\text{reward}}$$

+ Training Policy Gradients

1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward

$$W' = W - \nabla_{loss}$$

$$W' = W + \nabla \log P(a_t | S_t) R_t$$