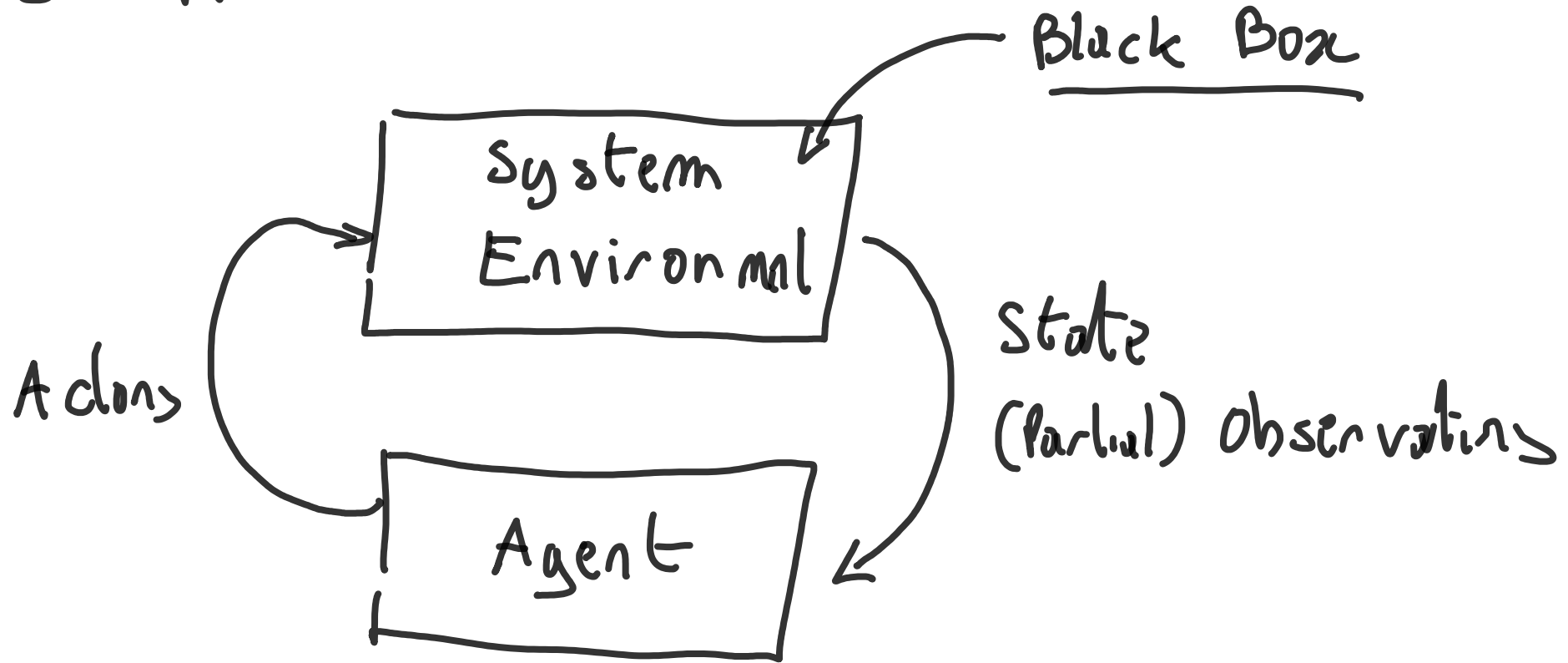


Reinforcement Learning for selective key applications in power systems



⇒ Goal: Finding a good way to pick actions (policy) from the interactions with the system

S_t : State of the world

A_t : Actions

For sake of simplicity, we assume that | the states are finite
| the actions are finite

$$P(S_{t+1} | A_t, S_t, A_{t-1}, S_{t-1}, \dots, A_0, S_0) \\ = P(S_{t+1} | A_t, S_t) \leftarrow \text{Markovian Assumption}$$

R_t : Rewards (instantaneous)

$$P(R_{t+1} | A_t, S_t, A_{t-1}, \dots) = P(R_{t+1} | A_t, S_t)$$

Policy: $\Pi: S_t, A_{t-1}, S_{t-1}, \dots, A_0, S_0 \mapsto A_t$

↳ Markovian Policy: $\Pi: S_t \rightarrow A_t$

Stochastic policy: $S_t \rightarrow \Pi(\cdot | S_t)$

Return (way of measuring the quality)

$\max_{\Pi} \mathbb{E}_{S_0} \mathbb{E}_{\Pi} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right]$ (Cumulative Reward)

↑
planning

discount factor $\gamma \in (0, 1)$

⇒ Everything is well defined

but far away future is less important than close future

Martovian Decision process (MDP) setting

$P(S_{t+1} | A_t, S_t)$ known

$$E_{\pi} \left(\sum \gamma^t R^t | s_0 \right) = E_{\pi}(R_1) + \gamma E_{\pi} \left[E_{\pi} \left[\sum_{t \geq 1} \gamma^{t-1} R^t | S_1, s_0 \right] \right]$$

value function

→ $V_{\pi}(s_0) = E_{\pi}(R_1) + \gamma E_{\pi} [V_{\pi}(s_1) | s_0]$

↳ V_{π} is the solution of a linear equation

Best Policy $Q^*(s,a) = r(s,a) + \sum_s P(s',s,a) \max_{a'} Q^*(s',a')$

Solution: → solving this system

Linear Programming

Iterative scheme

Fixed Point-Block Lemma

MDP

$P(S_{t+1} | S_t, A_t)$ is unknown

Model based approach: $P(S_{t+1} | S_t, A_t)$

Model free approach: No explicit estimation of $\hat{P}(S_{t+1} | S_t, A_t)$

↳ RL
Policy

Most classical
lot of proof

Value

Approximate optimal q^*

derive a policy

$$\pi^* = \underset{a}{\operatorname{argmax}} q^*(s, a) \leftarrow$$

Best policy
within a family

SOTA
Much
more modern