

(Discrete time Finite State)

Markov Decision Process (MDP)

↓ (choice) ↓ (Sequence in time)

Markov Assumption

{ For a Stochastic ^{discrete time} dynamical System }

$$P(X_{t+1} | X_t, U_t, X_{t-1}, U_{t-1}, \dots, X_0, U_0) = P(X_{t+1} | X_t, U_t)$$

"Given the present state, future ~~states~~ are independent of past ~~states~~"

"Conditioned on the present state future and past ~~states~~ are independent"

"State is the Complete Summary of the past"

⇒ It assumes that ^{for dynamical system} there exist a ~~state~~ but quantity (we are calling it State), such that ~~the~~ conditioned ~~on~~ the present state, future and past of the system is independent.

Markov Chain

→ S ∈ S state

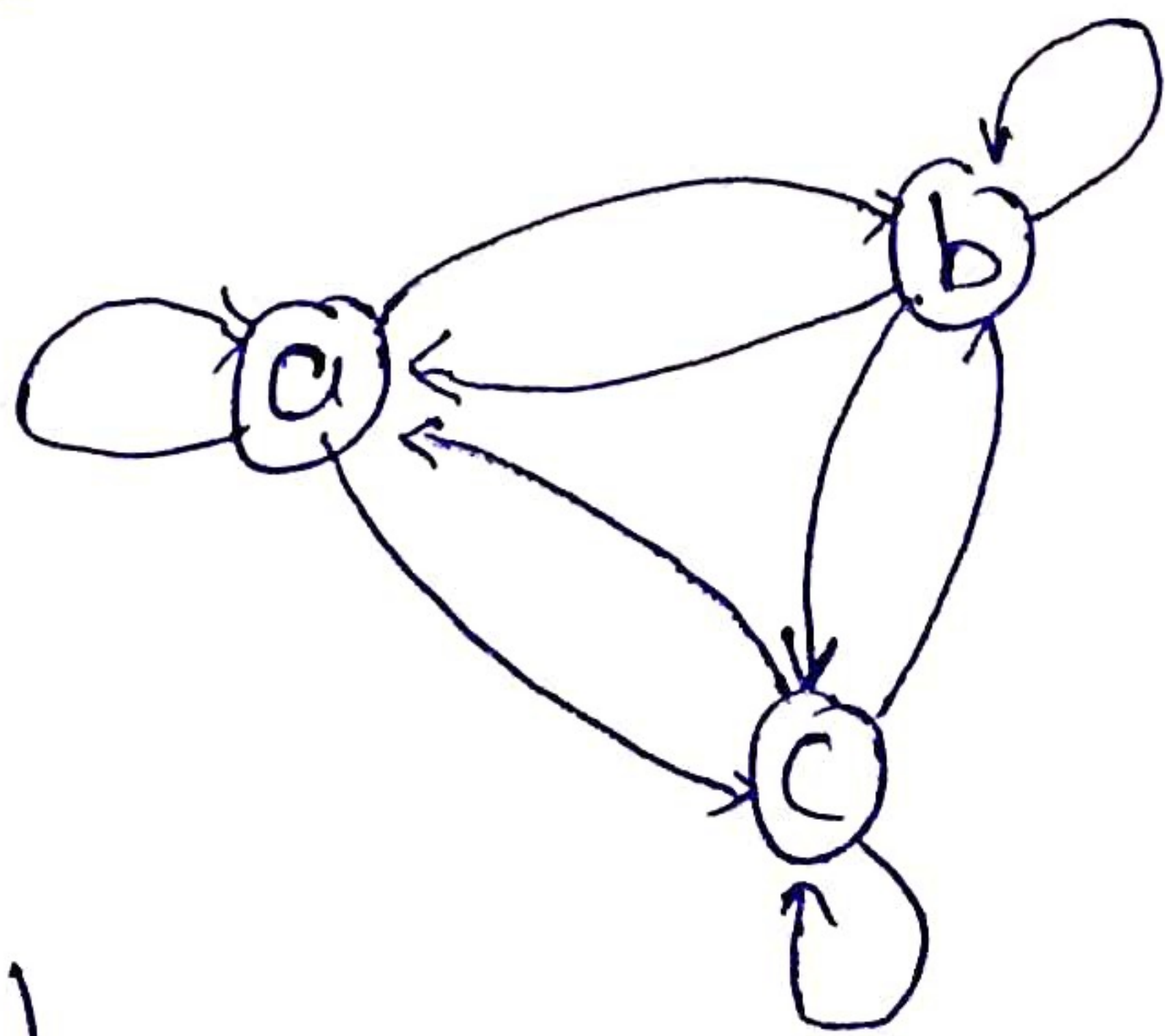
~~→ P(S)~~

→ $T(S, S') = P(S' | S)$

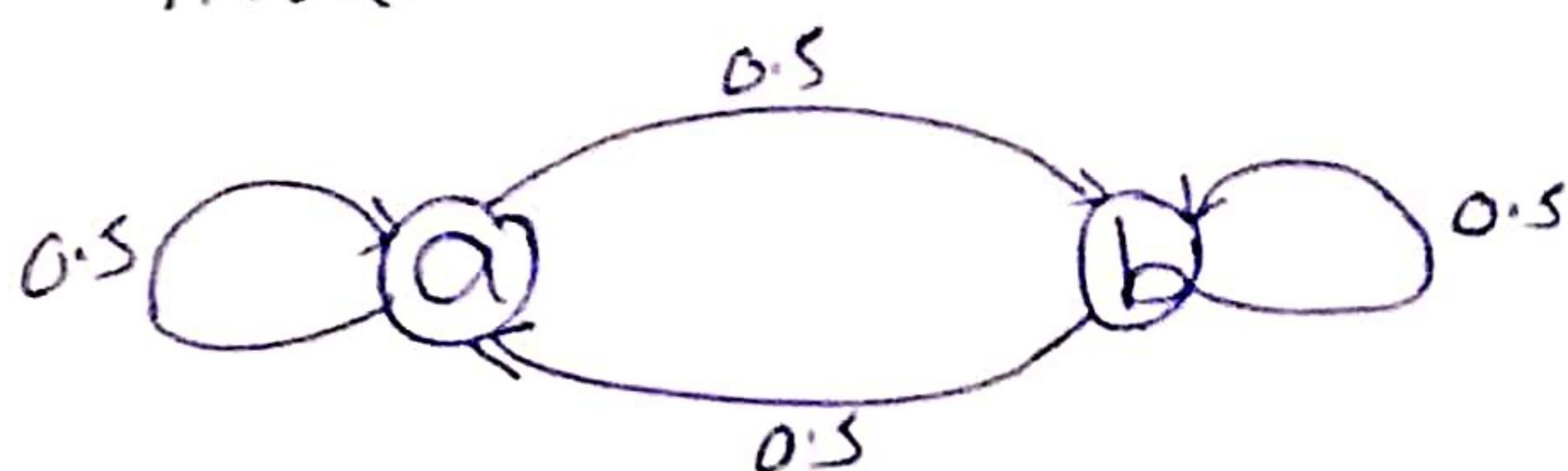
{ Probability that the system will transition to state S' , given its present state is S }

→ state state

in next time step



⇒ Markov chain is evolution of this kind of system over time.



a a b a b b a a b a b b

Sam Example of Markov Chain

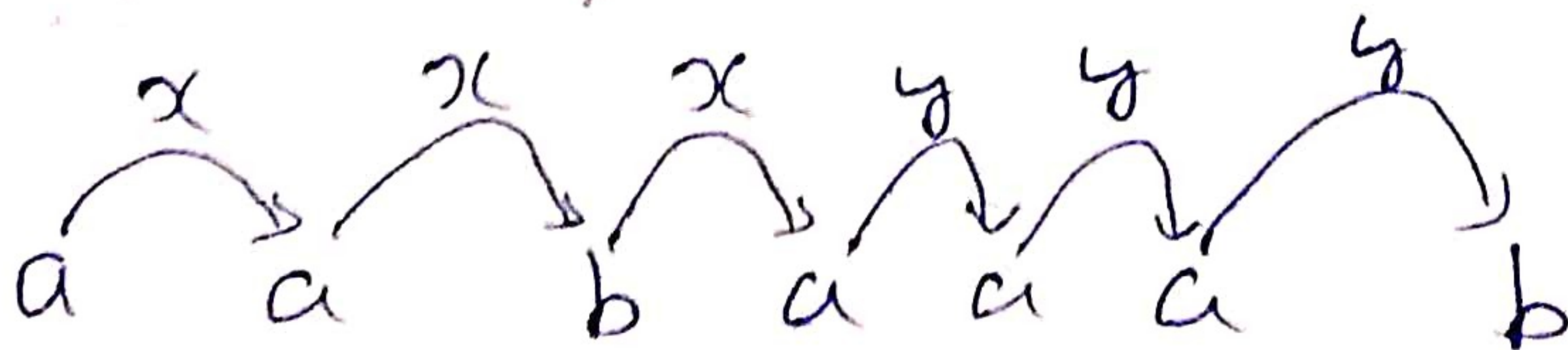
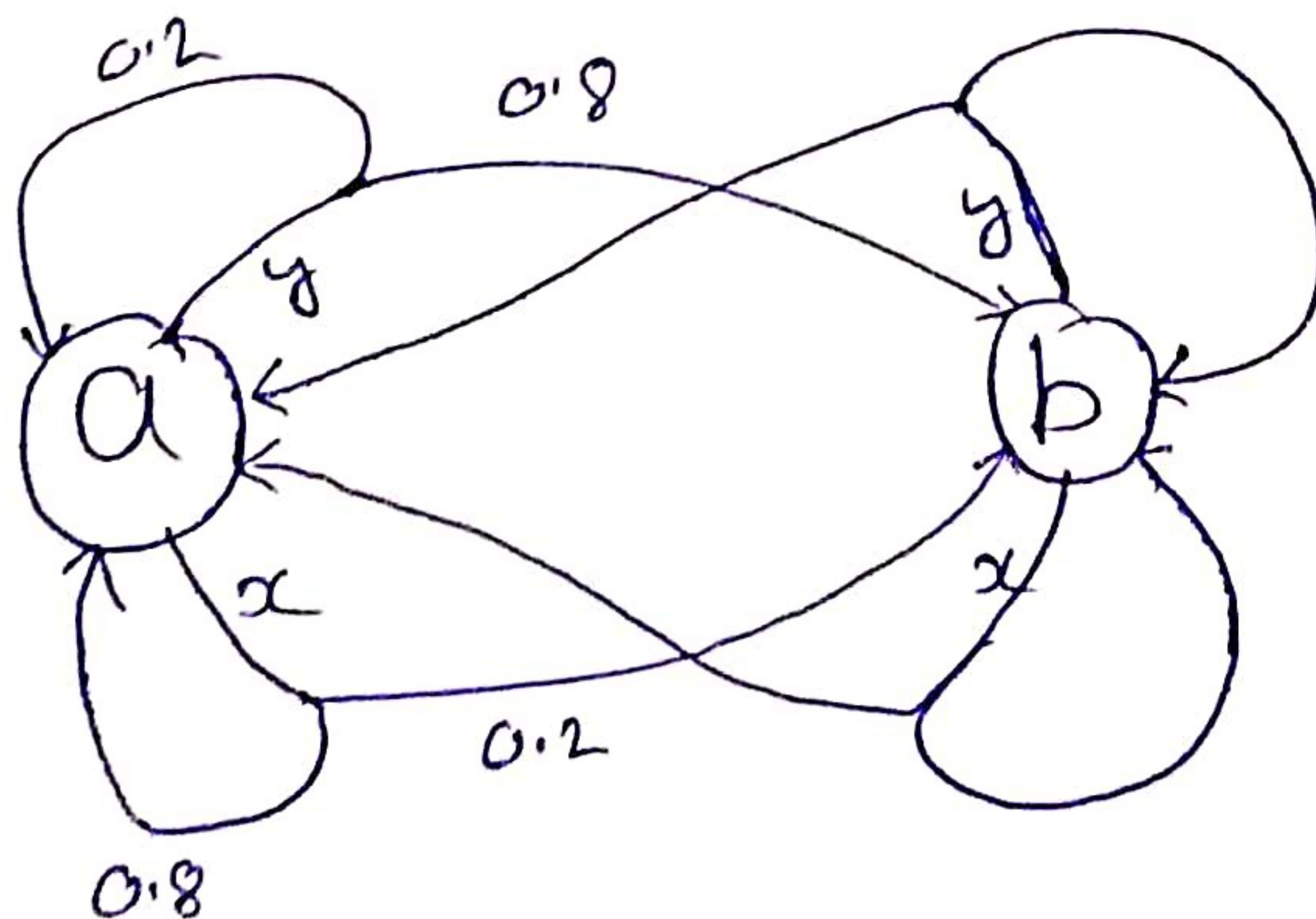
Markov Chain + Decision = Markov Decision Process

→ System with S States

→ a actions

→ $T(S, a, S') = P(S' | S, a)$

→ Start state

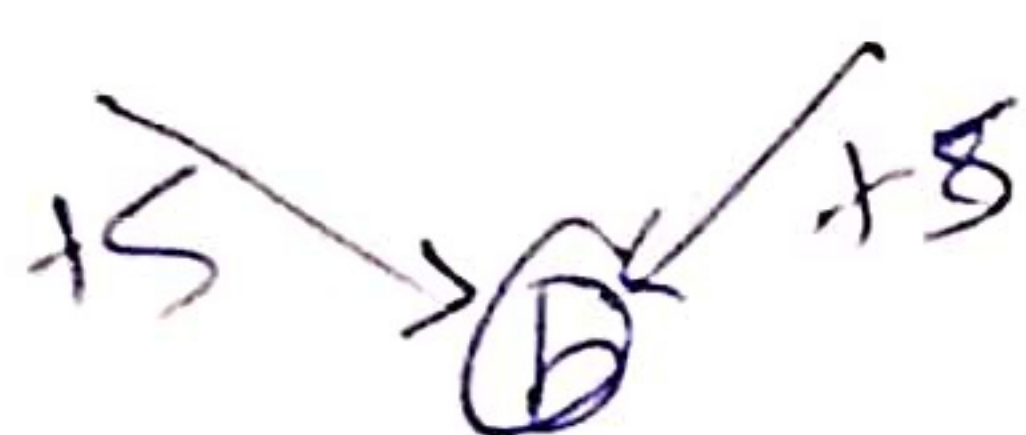


⇒ Evolution of this kind of system is referred to as Markov decision process.

⇒ You can also imagine ^{the system} ~~player~~ and ^{an} AI agent ~~making~~ making the choices/Decisions in every time step.

⇒ ~~To give a motivation to the agent~~ ^{To do that} we add rewards for transitions.

$R(s, a, s')$ { Sometimes, it only depends of state }
or sometimes the triple



⇒ Rational Agent : Select actions that maximize its (expected) utility.

{ Sum of ~~the~~ rewards }

→ It is also common to define utility as sum of discounted rewards.

↳ Idea : reward now is better than same reward in future.

$$r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

$$\boxed{\gamma \leq 1}$$

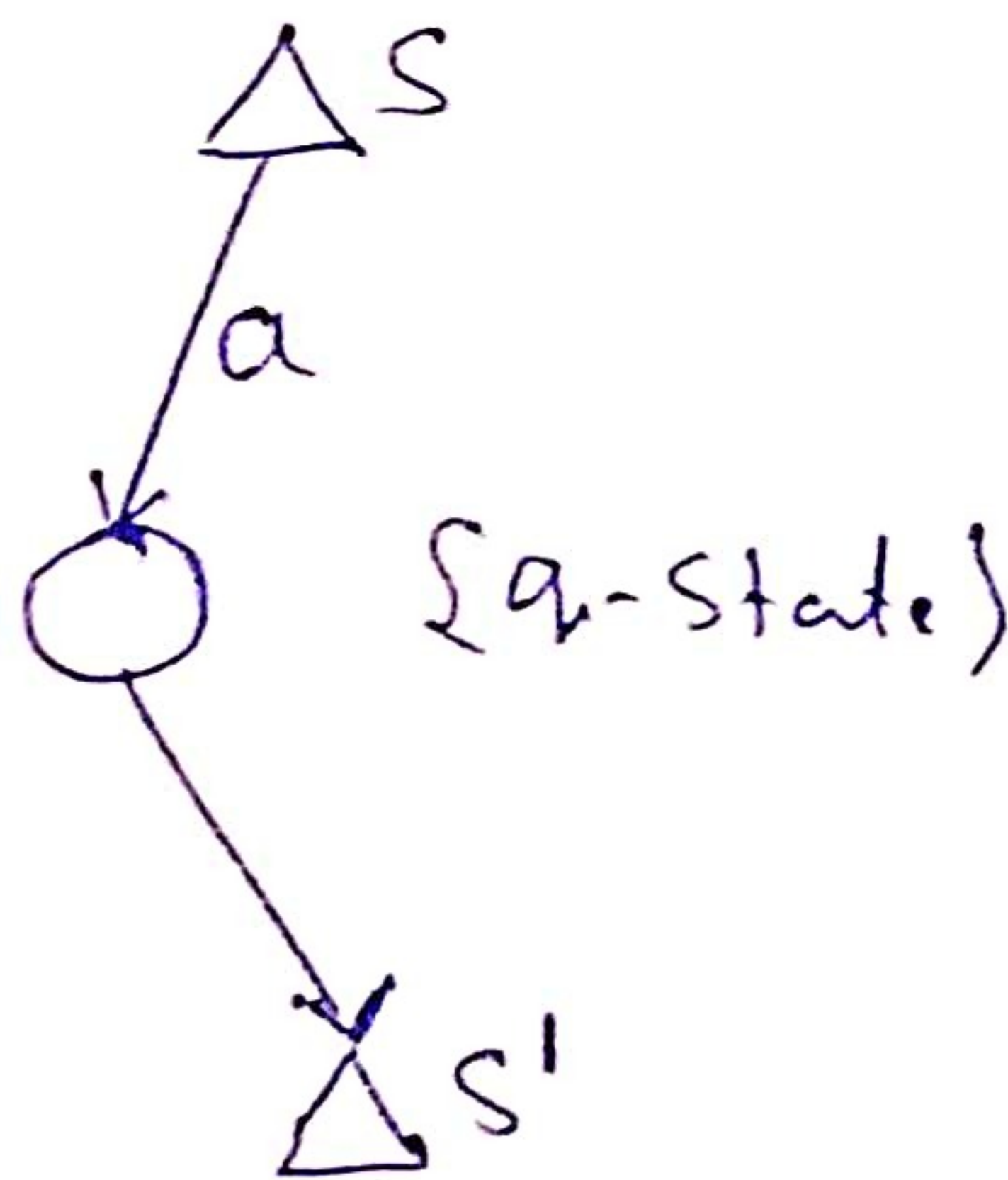
Policy

$$\pi: S \rightarrow A$$

(Optimal policy) π^* → Policy that maximizes the expected utility.

Plan \checkmark Policy

Sequence of Action



~~Value~~ Value of State

$V^*(S)$ → Expected utility starting at state S and acting optimally.

{ Similar we can define $Q^*(S, a)$ as the }
 { value of a-state }

$$V^*(s) = \max_a Q^*(s,a)$$

$$Q^*(s,a) = \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V^*(s')]$$

{ Bellman Equation }

Time Limited Value

$V_k(s)$ → Optimal value of state s if the game ends in k steps.

$$V_{k+1}(s) = \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V_k(s')]$$

* Value Iteration { An Algorithm for Solving MDP }

1. $V_0(s) = 0 \quad \forall s \in S$

2. Use Iterative Bellman Equation to compute V_1, V_2, \dots, V_n for all states.

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V_k(s')]$$

3. $\pi_{k+1}(s) \leftarrow \arg\max_a \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma V_k(s')]$

→ Policy at $k+1^{st}$ iteration.

⇒ Problem with value iteration is that, it is very slow!

- Complexity (as^2)
- actions (policy) converges way before value.

★ Policy Iteration

Let $V^\pi(s) \Rightarrow$ Expected utility starting at state s and following policy π .

1. Start with some initial policy π_0
2. Find value of all the states given the policy ~~and~~ by iteratively using the given below equation:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_{k+1}^{\pi_i}(s')]$$

3. Update the policy using the following.

$$\pi_{i+1} = \underset{a}{\operatorname{argmax}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

4. Repeat until policy converges

