# MDP : Part 2

**★ Convergance**

<u>Case 1:</u> If the three has maximum depth M, then $V_M$ holds the actual untruncted values.

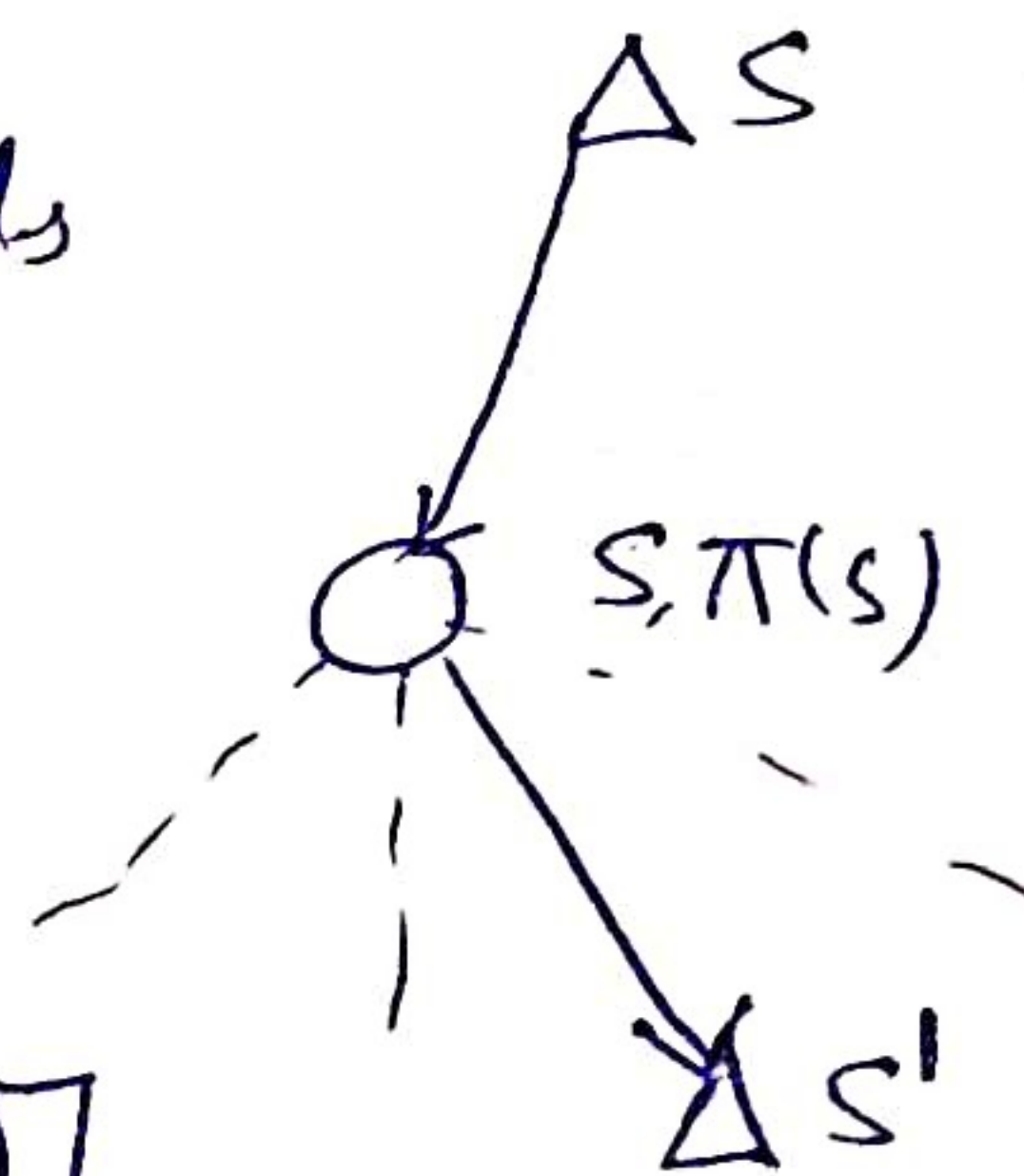<u>Case 2:</u> If the discount is less than 1

**★ Policy Evaluation**

⇒ Given the policy, you want to know, how good is that policy.

**★ Utilities for a fixed Policy**

$V^\pi(S)$ = Expected total discounted rewards starting in S and following $\pi$

$$V^\pi(s) = \sum_{s'} T(S, \pi(s), s')\left[R(S,\pi(s),s') + \gamma V^\pi(s')\right]$$

⇒ How do we calculate the V's for a fixed policy $\pi$?

  → <u>Idea 1:</u> Turn successive Bellman equations into updates.

  → <u>Idea 2:</u> Without the maxes, the bellman equations are just a linear system.

**★ Computing Actions from Values**

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

$$\underset{a}{argmax}$$

↳ Policy Extraction

**★ Computing Actions from Q-Values**

$$\pi^*(s) = \underset{a}{argmax}\, Q^*(s,a)$$

"actions are easier to select from" q-values than values!

**★ Problems with Value Iteration**

<u>Problem 1:</u> It's slow — $O(S^2A)$ per iteration

<u>Problem 2:</u> The "max" at each state rarely changes

<u>Problem 3:</u> The policy often converges long before the value.

**★ Policy Iteration**

<u>Step1:</u> (Policy Evaluation)

Calculate utilities for some fixed policy (not optimal utilities) until convergence

<u>Step2:</u> (Policy Improvement)

Update policy using one-step look-ahead with resulting converged (but not optimal) utilities as future values

$\Rightarrow$ Repeat steps until policy converges.

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

② $\pi_{i+1}(s) = \underset{a}{\text{argmax}} \sum_{s'} T(s, a, s')\left[R(s, a, s') + V^{\pi_i}(s')\right]$

$\Rightarrow$ Both Value Iteration and Policy Iteration are dynamic programs for Solving MDPs

**\* Reinforcement Learning**

$\Rightarrow$ Important ideas in reinforcement Learning:

* **Exploration**: You have to try unknown actions to get information

* **Exploitation**: Eventually, you have to use what you know

* **Regret**: Even if you learn intelligently, you make mistakes.

* **Sampling**: Because of chance, you have to try things repeatedly.

* **Difficulty**: Learning can be harder that Solving a known MDP.