

RL Week-4

LI DP : Value Iteration

- Turning the bellman optimality eqⁿ into an update rule.

$$v_{*}(s) = \max_a E[R_{t+1} + \gamma v_{*}(s_{t+1}) | S_t=s, A_t=a]$$

$$= \max_a \sum_{s', a'} p(s', a | s, a) [r + \gamma v_{*}(s')]$$

∴, write it as in terms of κ .

$$v_{\kappa+1}(s) = \max_a E[R_{t+1} + \gamma v_{\kappa}(s_{t+1}) | S_t=s, A_t=a]$$

$$= \max_a \sum_{s', a} p(s', a | s, a) [r + \gamma v_{\kappa}(s')]$$

$$v_0(s) = 0 \quad \forall s$$

$v_1(s)$ = is calculated $v_0(s)$

$v_2(s)$ = comes from $v_1(s)$

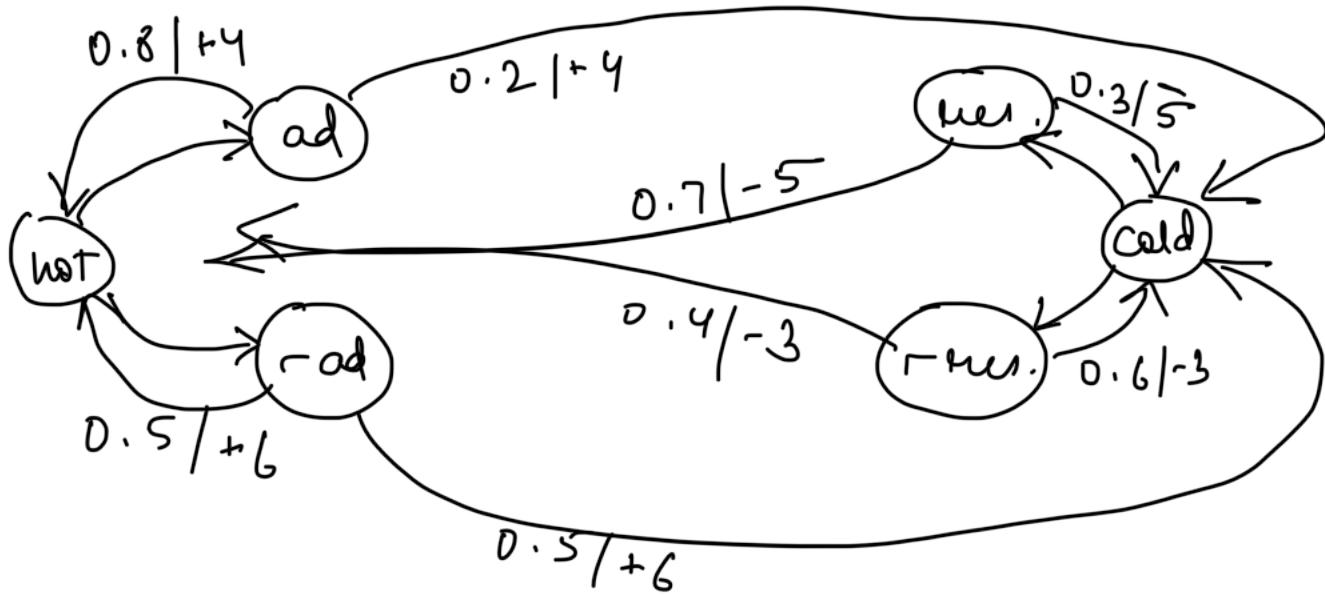
- a small threshold $\delta > 0$, that determines estimates accurately.

$$\text{Impt.} \rightarrow V(\text{terminal}) = 0$$

- output of a deterministic policy, $\pi^* \approx \pi^*$

$$\pi^*(s) = \arg \max_a \sum_{s', a} p(s', a | s, a) [r + \gamma V(s')]$$

- going back to manufacturer problem MDP.



$$v(\text{hot}) = v(\text{cold}) = 0 \quad [\text{Initialise}]$$

$$A(\text{hot}) = \{\text{ad}, \neg\text{ad}\}$$

$$A(\text{cold}) = \{\text{neu.}, \neg\text{neu.}\}$$

$$v(\text{hot}) = 6, \quad v(\text{cold}) = -3$$

(highest of all reward in both cases)

L2 Value Iteration

$$v_{k+1}(s) = v_{\pi}^{(1)}(s)$$

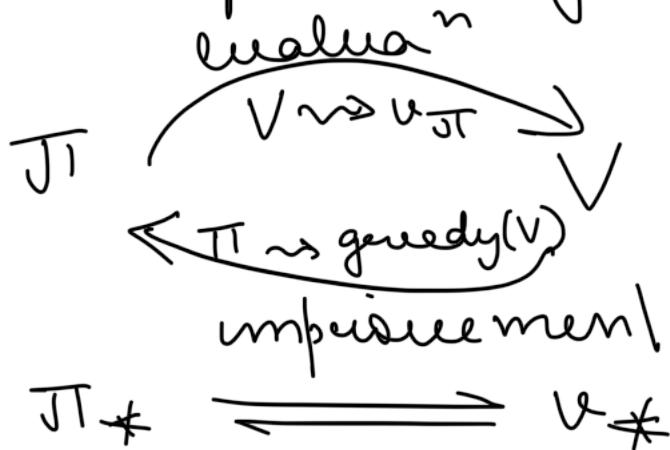
further demo, convergence happens after many iterations.

L3 More on DP

- asynchronous DP

1. Old algos. \rightarrow we had to do the update over entire state set. (Problem)

2. in async. DP \rightarrow the updates are not done for entire state set for each episode.
 3. in async. DP, we just have to sweep state \gg enough that it reaches convergence.
 4. in this way, order of the update can be chosen in a more flexible fashion.
 5. interestingly, it takes into account real time updates from envt. & DP.
 6. it helps you to update those parts which are specifically agent ~ space.
- if policy evaluaⁿ is stopped after 1 update , then we have to apply greedificaⁿ to that particular state . We get value iteraⁿ.
 - GPI (Generalized Policy Iteration)
 - \hookrightarrow in this , we allow policy evaluaⁿ & improvement interact with each other , independently .



- almost RL methods are kind of GPI.
- Policy iteration runs the update "forward" completion, after this improvement starts. In value iteration, 1 step of "update" is done before ~~the~~ improvement. But in GPI, asynch DP, the two methods work together.

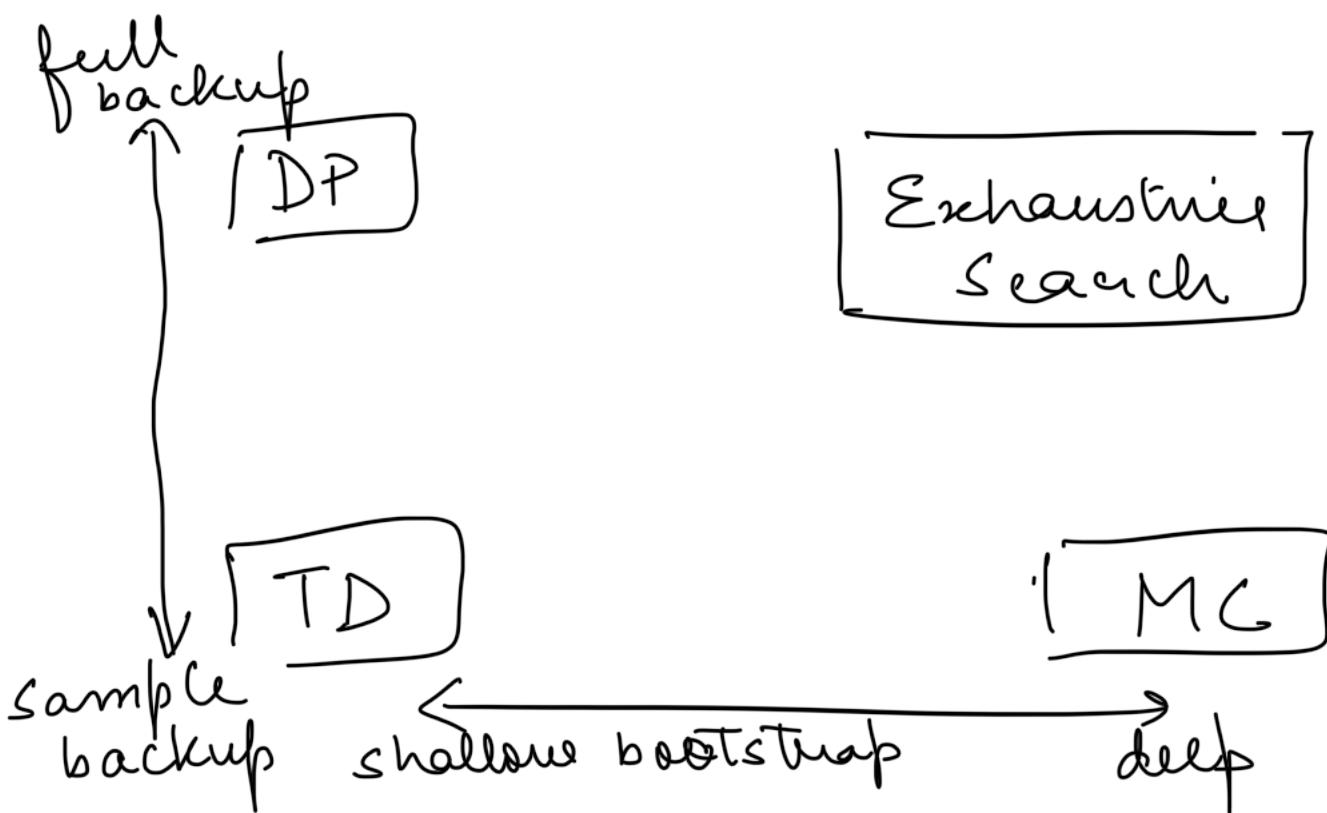
MC (Monte Carlo) Methods

- as per DP,
$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', a'} p(s', a' | s, a) [r + V v_{\pi}(s')]$$
- as per MC RL,
$$v_{\pi}(s) = E_{\pi} [G_t | S_t = s] \quad \forall s \in S$$
 1. learning directly from sample episodes of experience.
 2. not aug. (full mode) and is model free.
 3. no need of bootstrapping
 4. value func^w is calculated as the mean of discounted returns. (G_t)
- MC Periodic^{ms} -
 1. 1st visit $\rightarrow V(s)$ is the avg. of returns, first visit to s .

2. Ceeey visit \rightarrow returns the aug.
following all visits to s.

L5 Comparison of DP, MC & TD

- Bootstrapping
 - ↳ update using an estimate
- 1. DP bootstraps
- 2. MC does not bootstrap
- 3. TD bootstrap
- Sampling \rightarrow update calculated using samples w/o model.
 - 1. MC samples & TD
 - 2. DP does not sample.



Useful table -

	Bootstrap	Sampling
1. DP	✓	X
2. MC	X	✓
3. TD	✓	✓