

③ a) For an optimal policy  $\pi^*$  having cost  $J^*$

$$T_{\pi^*} J^* = T J^* \Rightarrow J^* = T J^* \quad (\text{Prop 2})$$

and  $J^*$  is the unique solution to the above equation ①

additionally,  $J^*$  is the unique solution to

$$T_{\pi^*} J^* = J \quad \text{--- ②}$$

Upon introducing a new action  $a'$  in a state  $s'$ , we get a new MDP. Let the optimal policy of such an MDP be  $\pi_{\text{new}}^*$  and cost be  $J_{\text{new}}^*$ .

From ②, since we don't have a change in transition probabilities &  $\pi^*(i)$  is also unchanged

$$T_{\pi^*} J^* = J^* \text{ still holds --- ③}$$

We still need to check if

$$J^* = T J^* \text{ to know whether } \pi^* \text{ is indeed } \pi_{\text{new}}^*.$$

[Assuming minimisation problem for part a)  
- although it really doesn't matter here]

$\therefore J^*$  was optimal for the older MDP

$(TJ^*)(i) = J^*(i)$  is older MDP

$$\Rightarrow J^*(i) = \min_a \sum_{+j} (P_{ij}(a) (g(i, a, j) + J^*(j)))$$

In the new MDP,

$$J_{\text{new}}^*(i) = \min_a \sum_{+j} (P_{ij}(a) (g(i, a, j) + J^*(j)))$$

holds for all  $i \neq s'$  ——— (4)

(because nothing has changed  
— no new actions have been added —)

claim (1):  $\nexists J^*(s') = \min_a \sum_{+j} (P_{s'j}(a) (g(s', a, j) + J^*(j)))$

$$= \min_a \sum_{+j} (P_{s'j}(a) (g(s', a, j) + J^*(j)))$$

then,  $J^* = J^*$  optimal &  $J^*(s') = J^*$  optimal

Proof: Procedure to test:

Assume that compare  $J^*(s')$  with

$$\sum_{+j} P_{s'j}(a') (g(s', a', j) + J^*(j))$$

if  $J^*(s') \geq \sum_{+j} P_{s'j}(a') (g(s', a', j) + J^*(j))$

then  $\pi^*(s')$  is still the solution to the eqn because substituting  $\pi^*(s')$  gives us a lower cost than  $q''$

$\Rightarrow \pi^*(s')$  is still optimal

$\therefore$  the claim is satisfied.

$$\Rightarrow \pi^* = \pi^*_{\text{new}} \text{ \& } J^* = J^*_{\text{new}}$$

Proof for the claim ①:

$$J^*(s) = \min_a \sum P_{s'j}(a) (g(s', a, j) + J^*(s'))$$

then combining with eqn ④, we have

(for new MDP),

$J^* = T J^* \Rightarrow J^*$  is still the optimal cost

$\therefore$  we already showed that expected cost of policy  $\pi^*$  is  $J^*$

$$\Rightarrow J^*_{\text{new}} = J^* \text{ \& } \pi^*_{\text{new}} = \pi^*$$

Summary: For the  $s'$  state equation, check if action  $a$  produced a minimum <sup>in RHS</sup>, if it doesn't the old policy remains the optimal policy.

claim 2: If  $a'$  gives a lower value in  
 $J^{a'} + J^*$  eqn (note that  $J^*$  is the  
 optimum in old MDP, it is not so in the new  
 MDP)  
 for  $s'$  state; then for that state  $s'$ ,  
 the action  $a'$  is optimal.

Proof:

Suppose  $a'$  is not the optimal action,  
 then it is as good as dealing with MDP  
 without that action

$\Rightarrow$  this is the old MDP

$\Rightarrow \pi^* \leq J^*$  are optimal in the  
 new MDP then

$$\Rightarrow J^*(s) = \min_a \sum_j P_{s'j} (g(s', a, j) + J^*(s))$$

has solution  $a = \pi^*(s)$  instead of  $a'$

this is a contradiction because  $J^*(s) \neq \sum_j \sum P_{s'j}(\cdot)$   
To find new policy  $\Rightarrow$  the claim is true.

From claim 1 we can infer that  $a'$  is the  
 optimal action at  $s'$  if  $\pi^* \neq \pi^*$  new



$$\text{Let } \pi_0(i) = \begin{cases} \pi^*(i) & i \neq s' \\ a' & i = s' \end{cases}$$

Using this policy we can perform policy iteration & find the new optimal policy

[note that  $\pi_0(i)$  is proper since

$$J^* \leq T_{\pi_0} J^* \quad (\text{if } a' \text{ is optimal})$$

as shown in that eqn in claim ①]

Since we already found one optimal action, policy iteration might converge faster

b) Assuming it is a maximisation problem

we are given that the cumulative reward

$J_{\pi^*}(s^*)$  decreases by  $\Delta$

$$J_{\pi^*}(s^*) = E(g(s^*, \pi^*(s^*), j) + J_{\pi^*}(j))$$

Since no other information is given, ~~and I~~  
 would like to interpret it as  $J_{\pi^*}(s^*)$   
 has changed because of a  $\Delta$  decrease in  
 average expected <sup>transition</sup> cost from  $s^*$

$$\text{that } J = E(g_{\text{new}}(s^*, \pi^*(s^*), j)) \\ = E(g_{\text{new}}(s^*, \pi^*(s^*), j))$$

Because of this change in transition cost  <sup>$-\Delta$</sup>  from  $s^*$ ,  
 there are 2 implications

i)  $J_{\pi^*}$  might not be solution to  ~~$J_{\pi^*} = T_{\pi^*} J_{\pi^*}$~~   
 $J = T_{\pi^*} J$   
 i.e. it may not be the expected cost of policy  $\pi^*$

ii)  $\pi^*$  may not be optimal policy (because  
 returns from  $s^*$  has reduced, intuitively  
 there is a possibility that a policy that  
 never has lesser number of transitions to  $s^*$   
 may be the  ~~$\pi^*$~~  new optimal policy,  $\pi_{\text{new}}^*$ )

## Checking if $\pi^*$ is optimal

Per Step ①: Perform policy evaluation

$$J(i) = \sum_j P(j|\pi_k(i)) (g(i, \pi_k(i), j) + J(j))$$

$$\text{where } g(i, \pi_k(i), j) = \begin{cases} g_{\text{original}}(i, \pi_k(i), j) \\ g_{\text{original}}(s^+, a^+, j) - \end{cases}$$

$$J(i) = E(g_{\text{original}}(i, \pi_k(i), j) + J(j)) \quad \# i \neq s^+$$

there  $\epsilon$  for  $i = s^+$

$$J(s^+) = E(g_{\text{original}}(s^+, \pi_k(s^+), j) + J(s^+)) \quad \text{--- } \Delta$$

where  $g_{\text{original}}$  is the transition cost of

original MDP. Let soln to these eqs be  $J_{\pi^*}$

Step ②: Perform policy improvement:

$$T_{\pi_{\text{new}}} J_{\pi^*} = T J_{\pi^*}$$

$$\Rightarrow (T_{\pi_{\text{new}}} J_{\pi^*})_i = \max_a E(g(i, a, j) + J_{\pi^*}(j))$$

If  $J_{\pi_{\text{new}}}$

Get  $\pi_{\text{new}}$  & evaluate  $J_{\pi_{\text{new}}}$ .

$$\text{if } J_{\pi_{\text{new}}} = J_{\pi^*}$$

then  $\pi^*$  is the optimal policy

If  $J_{\pi_{\text{new}}} < J_{\pi^+}$ ,  $\pi^+$  is not optimal.

In order to ~~confirm~~ find the new optimal

policy we can do policy iteration

with  $\pi_{\text{new}}$  as the initial proper policy.

The expected cost will also be obtained via  
policy iteration in policy evaluation ~~step~~ <sup>step</sup>