

Ex 1 :

- .) 1. We can take  $\forall s \neq 1, s \neq 14 \quad r(s,a) = -\varepsilon$ .  
 with  $\varepsilon > 0$

$\varepsilon$  should be positive in order to avoid the infinite cycle policy, but it should not be too big, because the agent would then suicide directly. But by taking  $\varepsilon = 1$ , because the distance between the cases 1 and 14 is only of 4, the agent has no need to suicide if he appears on the case 2.

$$V(2) = \frac{-\varepsilon}{2} + \frac{-\varepsilon}{6} + \frac{-\varepsilon}{10} + \frac{10}{14} = -3\varepsilon + 10$$

$$\text{So } -3\varepsilon + 10 > -\varepsilon + V(1) = -10 \cdot \varepsilon \Rightarrow \varepsilon < \frac{20}{2} = 10$$

It is the same reasoning for  $V(5)$

So by taking  $\varepsilon = 1$  there is no risk to suicide

The resulting value is :

5	-10	7	6	
6	7	8	5	
7	8	9	4	
8	9	10	3	

↑ final state

2. We suppose a general MDP with reward  $r(s, a)$ .

If we take an affine transformation  $r'(s, a) = A(r(s, a)) + B$

The new value function is  $V' = AV + B$ . (because we keep the same policy). But the optimal policy is not the same.

Even if  $A > 0$ , by selecting the appropriate  $B$  we can change the policy. For example if  $A=1, B=-1000$  if we take the same environment as in Q1, the optimal policy is the quickest suicide.

3. Because  $r^>(s, a) = r(s, a) + 5 = 4 > \forall s \neq (1, 14)$

So the goal become equivalent to staying for eternity on the board by jumping the infinite cycle

$5 \rightarrow 6 \rightarrow 10 \rightarrow 9 \rightarrow 5$ . And the value is  $+\infty$

almost everywhere but in states 1 and 14.

D<sub>RL</sub>

$$\text{Eq2 } V^*(s) - V^{\pi_Q}(s) = \underbrace{Q^*(s, \pi^*(s)) - Q^*(s, \pi_Q(s))}_{A} + \underbrace{Q^*(s, \pi_Q(s)) - Q^{\pi_Q}(s, \pi_Q(s))}_{B}$$

$$A = Q^*(s, \pi^*(s)) - \underbrace{Q(s, \pi^*(s)) + Q(s, \underline{\pi^*(s)})}_{\text{O}} - Q^*(s, \pi_Q(s))$$

$$\leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s)) + Q(s, \underline{\pi_Q(s)}) - Q^*(s, \pi_Q(s))$$

Car  $\pi_Q(s) = \arg \max_a Q(s, a)$

$$\leq \max_s |Q^*(s, \pi^*(s)) - Q(s, \pi^*(s))| + \max_s |Q^*(s, \pi_Q(s)) - Q(s, \pi_Q(s))|$$

$$\leq 2 \max_{s, a} |Q^*(s, a) - Q(s, a)| \leq 2 \|Q^* - Q\|_\infty$$

$$\begin{aligned} B &= Q^*(s, \pi_Q(s)) - Q^{\pi_Q}(s, \pi_Q(s)) \\ &= \cancel{r(s, \pi_Q(s))} + \gamma \sum_{s'} P(s' | s, \pi_Q(s)) V^*(s') - \cancel{r(s, \pi_Q(s))} - \cancel{\sum P(s' | s, \pi_Q(s))} \\ &\leq \gamma \sum_{s'} P(s' | s, \pi_Q(s)) \|V^* - V^{\pi_Q}\|_\infty \\ &\leq \gamma \|V^* - V^{\pi_Q}\|_\infty \end{aligned}$$

$$\begin{aligned} \text{Done } \forall s, \quad V^*(s) - V^{\pi_Q}(s) &\leq 2 \|Q^* - Q\|_\infty + \gamma \|V^* - V^{\pi_Q}\|_\infty \\ \Rightarrow \|V^* - V^{\pi_Q}\| &\leq 2 \|Q^* - Q\|_\infty + \gamma \|V^* - V^{\pi_Q}\|_\infty \\ \Rightarrow \|V^* - V^{\pi_Q}\| &\leq \frac{2}{1-\gamma} \|Q^* - Q\|_\infty \quad (\gamma < 1) \\ \Rightarrow \forall s, \quad V^*(s) - V^{\pi_Q}(s) &\leq \frac{2}{1-\gamma} \|Q^* - Q\|_\infty \end{aligned}$$

$$\text{et } \boxed{V^{\pi_Q}(s) \geq V^*(s) - \frac{2}{1-\gamma} \|Q^* - Q\|_\infty}$$

Q3 We develop the wanted expression

$$\begin{aligned} \textcircled{1}) & \sum_s \mu^{\pi'}(s) \sum_a (\pi'(a|s) - \pi(a|s)) Q^{\pi'}(s, a) \\ \textcircled{2}) & = \sum_s \mu^{\pi'}(s) \sum_a (\pi'(a|s) - \pi(a|s)) \left( r(s, a) - g^{\pi'} + \underbrace{\sum_{s', a'} p(s'|s, a) \pi(a'|s') Q^{\pi'}(s', a')}_{B} \right) \\ & = A + B \end{aligned}$$

$$\begin{aligned} A &= \sum_{s, a} \mu^{\pi'}(s) (\pi'(a|s) - \pi(a|s)) (r(s, a) - g^{\pi'}) \\ \textcircled{3}) &= \sum_{s, a} \mu^{\pi'}(s) \left( \underbrace{\pi'(a|s) r(s, a)}_{\text{1}}, \underbrace{- \pi'(a|s) g^{\pi'}}_{\text{2}} - \pi(a|s) r(s, a) + \pi(a|s) g^{\pi'} \right) \end{aligned}$$

We remark

$$\text{that } \sum_{s, a} \mu^{\pi'}(s) \pi'(a|s) r(s, a) = g^{\pi'}$$

We sum first on a  
and then on s  $\rightarrow g^{\pi'} - g^{\pi} = 0$

$$\textcircled{4}) = g^{\pi'} - \sum_{s, a} \mu^{\pi'}(s) \pi(a|s) r(s, a)$$

$$\begin{aligned} B &= \sum_{\substack{s, a \\ s', a'}} \mu^{\pi'}(s) (\pi'(a|s) - \pi(a|s)) \underbrace{p(s'|s, a) \pi(a'|s') Q^{\pi'}(s', a')}_{\text{const of } s, a}, \end{aligned}$$

$$\begin{aligned} \textcircled{5}) &= \sum_{s', a'} \pi(a'|s') Q^{\pi'}(s', a') \underbrace{\sum_{s, a} \mu^{\pi'}(s) (\pi'(a|s) - \pi(a|s)) p(s'|s, a)}_{\mu^{\pi'}(s')} - \sum_{s, a} \mu^{\pi'}(s) \pi(a|s) p(s'|s, a) \end{aligned}$$

$$\begin{aligned} &= - \sum_{s, a} \mu^{\pi'}(s) \pi(a|s) \sum_{s', a'} \underbrace{\pi(a'|s') Q^{\pi'}(s', a') p(s'|s, a)}_{\text{Bellman}} \end{aligned}$$

$$+ \sum_{s', a'} \pi(a'|s') Q^{\pi'}(s', a') \mu^{\pi'}(s')$$

$$B = - \sum_{s,a} \mu^{\pi'}(s) \pi(a|s) Q^{\pi}(s,a) - \sum_{s,a} \mu^{\pi'}(s) \pi(a|s) (-r(s,a) + g^{\pi})$$

$$+ \sum_{s',a'} Q^{\pi}(s',a') \pi(a'|s') \mu^{\pi'}(s')$$

$$= - \sum_{s,a} \mu^{\pi'}(s) \pi(a|s) (g^{\pi} - r(s,a))$$

$$= -g^{\pi'} + \sum_{s,a} \mu^{\pi'}(s) \pi(a|s) r(s,a)$$

$$\boxed{A+B = g^{\pi'} - g^{\pi}}$$

Ex 5 We have a 6 story building with 2 elevators.

(i) We have to define the tuple  $\Pi = (S, A, p, r)$ .

(ii) 1.  $S = (z_t^1, z_t^2)$  with no history,  
 $(z_t^1, z_t^2) \in [0, 6]^2$ .

We now do not modelise the time the elevators takes to go from one story to another one. And we use a discrete time. One unit of time is the time corresponding to opening the doors, entering and exiting the elevator, pushing the button, waiting to close the door and going to another story, not necessarily the one just above or just below.

This modellisation is valid if the elevator is fast (to go from a story to another one) but the doors and the people are slow.

But we have also to modelise the different people.

$S = ((z_t^1, z_t^2), (W, E))$ .  $W_{i,j}^t$  is the number of people currently waiting on the  $i$ -th story willing to go to the  $j$ -th story.  $E_{i,t}^e$  is the number of people currently in the elevator  $e$  (which is currently in the story  $z_t^e$ ) willing to go to the story  $i \neq z_t^e$ .

The set of possible actions is to go to another story.

So  $A = ((g_{t+1}, b_{t+1}) \in [0, 6]^2)$

If the elevator was slow or the number of story was much bigger, we would choose

$A = ((d_{z_1}, d_{z_2}) \in [-1, 1]^2)$ , but that is not the case in our assumptions.

$p(s'|s, a)$ ? We consider the elevator is responding perfectly to the orders of the agents. So  $z_{t+1}^e = a_t^e$  with no stochasticity. We make the hypothesis the elevators can contain any number of people. And people in the story just want to go on the floor, and they enter the first elevator they see. They do not exit an elevator unless they are at the right story.

$$\text{So } E_{z,t+1}^e = E_{z,t}^e + w_{z,t,z}^t \quad \text{if } s \neq z_{t+1}^e$$

"the number of people in elevator  $e$ , willing to go to story  $z$  equals the previous number of people willing to go to  $z$  plus the people willing to go to  $z$  who were on the precedent story."

$$E_{z,t+1}^e = 0 \quad \text{if } z = z_{t+1}^e$$

"They exited the elevator"

$$w_{z_1, z_2}^{t+1} = w_{z_1, z_2}^t + N_{z_1, z_2}^t \quad \text{if } z_{t+1}^e \neq z_1$$

"People who were waiting are still waiting and did not take the stairs if the elevator did not come"

$$w_{z_1, z_2}^{t+1} = N_{z_1, z_2}^t \quad \text{if } z_{t+1}^e = z_1$$

"People who were waiting took the elevator and a bunch of new people arrived"

\* NB = if elevator 1 and 2 simultaneously come to a story) the waiting people all choose the elevator with the smallest number of people within.

$N_{i,j}^t$  is the number of people coming between time  $t$  and  $t+1$ .

So they always wait until  $t+1$ , and if the elevators come, they enter in.

$N_{i,j}^t$  follows a Poisson law  $P(\lambda_{i,j}^t)$

If we want the process to be Markovian, we should not consider the time dependence of  $\lambda_{i,j}^t \rightarrow \lambda_{i,j} = c_i$ .

So in our modelisation, there is no Night and Day cycle.

Finally, the reward function is  $r(s, a, s')$ : we can take a penalty function of the number of people currently waiting.

$$r(s_t, a, s_{t+1}) = - \sum_{i,j} w_{i,j}^t - \sum_i E_{i,t}^e \quad \forall a$$

The penalty is linear in the total waiting time.

If we choose instead ~~the~~ to modelise the penalty as the sum of the squares of the individuals waiting time, we could not do it unless we modelise each individual.

Ex5 1 : See Code

2. See Code

3. For proba-succ = 0,8 : (stochastic)

The value iteration terminated with 1194 iterations.

The policy iteration terminated in 9 iterations.

And the value function converges to 691.

For proba-succ = 0,99 (almost non stochastic)

The value iteration terminated with 1175 iterations.

The policy iteration terminated with 12 iterations.

And the value function converges to 545.

545 is smaller than 691 which is normal, because

the non stochastic environment is less risky.

the number of iterations in the value iteration is almost the same

But the number of iterations decreases in the policy iteration.  
if there is stochasticity.

If we are in the non-stochastic environment, there is ... multiplicity of the optimal policy, because it is equivalent to go first on the right and then in the bottom. or to go on the bottom and then on the right.

But in the stochastic environment,  $\rightarrow \downarrow$  or  $\uparrow$  are not the same because cases containing a trap.

So the optimal policy in the stochastic environment tries to reduce the time on the frontier of the cases with a trap.

h. Value Iteration needs more iterations but each iteration is quick to run. Policy iteration needs less iterations but each iteration takes more time

~~There is no unicity in the policy iteration, but there is unicity in the value iteration.~~

The time needed to compute the value iteration is longer (22s) than the time needed to compute policy iteration (0,3s).

The stopping criterion in the policy iteration is straightforward (equality of the policies) whereas the stopping criterion of the value iteration is a little bit more technical, because only asymptotic.