

A11) We can write

$$p(\pi' | s, a) = \sum_{\pi, s', a', s''} p(\pi', s'' | s', a').$$

$$\pi(a' | s') \cdot p(\pi, s' | s, a)$$

which can be simplified as :

$$p(\pi' | s, a) = \sum_{a', s'} p(s' | s, a) \pi(a' | s') p(\pi' | s', a')$$

Hence  $R_{t+2}$  is dependent on  $S_t, A_t$

~~Although we can also write  $R_{t+2}$  as~~

$$A12) E[R_{t+2} | s, a] =$$

$$\sum_{\pi'} \pi' \sum_{\pi, s', a', s''} p(\pi', s'' | s', a') \pi(a' | s') p(\pi, s' | s, a)$$

$$\text{A13) } V_{\pi}(s) = E[G_t | S_t = s]$$

$$V_{\pi}(s) = E[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \sum_a E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$\pi(a|s)$$

$$= \sum_{a, \pi, s'} E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r]$$

$$\pi(a|s) p(\pi, s' | s, a)$$

$$= \sum_{a, \pi, s'} \left( E[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] + \gamma E[G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] \right)$$

$$\cdot \pi(a|s) p(\pi, s' | s, a)$$

$$V_{\pi}(s) = \sum_{a, \pi, s'} \pi(a|s) p(\pi, s' | s, a) \cdot \{ r + \gamma V_{\pi}(s') \}$$

AA 14)

$t$	1	2	3	4
$R_t$	2	-1	10	-3
$G_t$	3.625	3.25	8.5	-3

We know that

$$G_t = \sum_{T=0}^{\infty} \gamma^T R_{t+T+1}$$

Now  $R_t = c \quad \forall t$

$$G_t = c \sum_{T=0}^{\infty} \gamma^T$$

$$G_t = \frac{c}{1-\gamma}$$

A15) We will choose the action that maximizes the expected return.

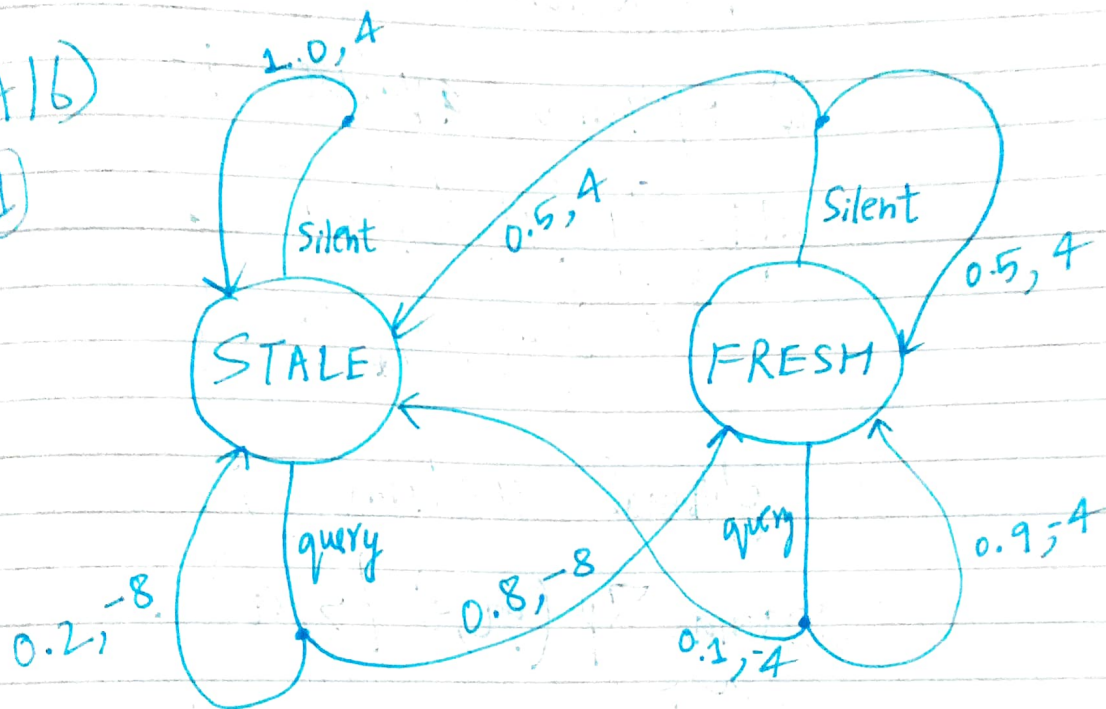


$$\pi_*(s) = \operatorname{argmax}_a E[G_t | S_t = s, A_t = a]$$

$$= \operatorname{argmax}_a \sum_{r, s'} p(r, s' | s, a) \{ r + V_*(s') \}$$

A16)

①



②

$$\gamma = 0.5$$

$$\text{reward\_stale} = -10.0$$

$$\text{reward\_fresh} = +10.0$$

$$T = 3$$

value( $a^*$ )	$t=0$	$t=1$	$t=2$	$t=3$
stale	5.75(s)	3.5(s)	-1.0(s)	-10.0
fresh	6.3(s)	4.75(s)	4.0(s)	10.0

Hence optimal policy is to stay  
Silent always and never query.

Steps: i) Use dynamic programming,  
start calculation with  $t=3$   
backwards

ii) Optimal value is

$$\max_a \sum_{r,s'} p(r,s' | s,a) \cdot \{r + V_*(s')\}$$

Optimal action is

$$\operatorname{argmax}_a \sum_{r,s'} p(r,s' | s,a) \cdot \{r + V_*(s')\}$$

③ Value iteration :

k	v(stale)	v(fresh)
0	0.0	0.0
1	4.0	4.0
2	6.0	6.0
3	7.0	7.0
4	7.5	7.5

Using 
$$V_{k+1}(s) = \max_a \sum_{\pi, s'} p(\pi, s' / s, a) \cdot \{ \pi + \gamma V_k(s') \}$$

Policy iteration

Initially,  $v(\text{stale}) = v(\text{fresh}) = 0$

$\pi(\text{stale}) = \text{silent}$

$\pi(\text{fresh}) = \text{silent}$

iteration 1

evaluation :

$$v(\text{stale}) = 4 + 0.5 v(\text{stale})$$

$$\Rightarrow \boxed{v(\text{stale}) = 8}$$



$$v(\text{fresh}) = 0.5 \left( 4 + 0.5 v(\text{stale}) \right) + 0.5 \left( 4 + 0.5 v(\text{fresh}) \right)$$

$$v(\text{fresh}) = 8$$

policy ~~iteration~~ improvement:

$$\pi(\text{fresh}) = \text{silent}$$

$$\pi(\text{stale}) = \text{silent}$$

policy is stable so all iterations ahead are same.

Used Bellman Equations for evaluation  
and  $\pi(s) = \operatorname{argmax}_a \sum_{\pi, s'} p(\pi, s' / s, a) \cdot$

$$\{ \pi + \gamma v(s') \}$$

for improvement