

Tutorial -1 :

1. $P_t(s'|s)$, $r_t(s)$ if system occupies state s at time t
 either don't do anything = ϕ
 or choose from A_s
 s.t. $P_t(s'|s,a)$ reward $r_t(s,a)$

(a) $(T, \mathcal{S}, A_s, s \in \mathcal{S}, r'_t(s,a), P'_t(j|s,a))$ is MDP
 true! let $T = \tau$ (given discrete-time MDP) ($N = \infty$ or $N < \infty$)
 (i.e. $T = \{1, 2, \dots, N\}$)

now $\mathcal{S}' = \text{given set of all states } S$

now $A'_s = A_s \cup \{\phi\}$, as A_s given and we have one more action $\forall s \in \mathcal{S}'$
 ϕ s.t. we don't do anything

now, $r'_t(s,a) = \begin{cases} r_t(s) ; a = \phi \\ r_t(s,a) ; a \neq \phi \end{cases}$ this is given

$$P'_t(j|s,a) = \begin{cases} P_t(j|s) ; a = \phi \\ P_t(j|s,a) ; a \neq \phi \end{cases}$$

(b) $T=2$, $v(j) = \text{terminal reward}$ true

$$\text{MTR}(s_1) = \max_{a' \in A_s \cup \{\phi\}} (r'_t(s_1, a') + \sum P'_t(j|s_1, a') v(j))$$

now this is for given s_1

$$\text{now, } \text{MTR}(s_1) = \max_{a' \in A_s} \left\{ r_t(s_1, a') + \sum_{\substack{j \\ a' \in A_s}} P_t(j|s_1, a') v(j), r_t(s_1) + \sum_{\substack{j \\ a' \in A_s}} P_t(j|s_1) v(j) \right\}$$

case $a' \in A_s$ case ϕ under

we try to pick $a' \in A_s \cup \{\phi\}$ s.t. it maximizes the above
 so let:

$$a^*(s) = \operatorname{argmax}_{a' \in A_s \cup \{\phi\}} \left\{ r'_t(s, a') + \sum P'_t(j|s, a') v(j) \right\}$$

true $a^*(s)$ is function s.t. for state $s \in \mathcal{S}$ what
 action to do at $T=1$ to maximize reward

2. $T = \{1\}$

$$\begin{aligned} \mathcal{S} &= \{s_1, s_2\} & A_s &= \{a_{11}, a_{12}\} \\ & & A_{s_2} &= \{a_{21}, a_{22}\} \\ r_1(s_1, a_{11}) &= 5 & r_1(s_2, a_{21}) &= -1 \\ r_1(s_1, a_{12}) &= 10 & r_1(s_2, a_{22}) &= \end{aligned}$$

	s_1	s_2
(s_1, a_{11})	0.5	0.5
(s_1, a_{12})	0	1
(s_2, a_{21})	0.8	0.2
(s_2, a_{22})	0.1	0.9

(a) $v(s) = 0 \quad \forall s \in \mathcal{S}$ true

$$a^* = \operatorname{argmax}_{a' \in A_s} \left\{ r_1(s, a) + \sum_{j \in \mathcal{S}} v(j) P_1(j|s, a) \right\} \xrightarrow{0}$$

$$\Rightarrow a^* = \underset{a \in A_S}{\operatorname{argmax}} \gamma_1(s, a)$$

now if at s_1 at $t=1$ then

$$a^*(s_1) = \underset{a \in \{a_{11}, a_{12}\}}{\operatorname{argmax}} \gamma_1(s_1, a)$$

$$a^*(s_1) = a_{12} \text{ as } \gamma_1(s_1, a_{12}) > \gamma_1(s_1, a) \forall a \in A_S,$$

$$\text{now } a^*(s_2) = \underset{a \in \{a_{21}, a_{22}\}}{\operatorname{argmax}} \gamma_1(s_2, a) = a_{22}$$

$$\text{so, } a^* = \begin{cases} a_{12}; \text{ at } t=1 \text{ at } s_1 \\ a_{22}; \text{ at } t=1 \text{ at } s_2 \end{cases}$$

$$(b) \text{ now total reward} = \gamma_1(s, a') + \sum v(s') p_i(j|s, a')$$

$$\text{now, if at } s_1: \text{TR} = \gamma_1(s, a') + v(s_1) p_i(s_1|s_1, a) + v(s_2) p_i(s_2|s_1, a)$$

$$\begin{aligned} a^*(s_1) &= \underset{a \in \{a_{11}, a_{12}\}}{\operatorname{argmax}} \left(5 + d \times 0.5 + e \times 0.5, 10 + d \times 0 + e \times 1 \right) \\ &= \underset{a \in \{a_{11}, a_{12}\}}{\operatorname{argmax}} \left(5 + \frac{d+e}{2}, 10+e \right) \end{aligned}$$

$$\begin{aligned} a^*(s_2) &= \underset{a \in \{a_{21}, a_{22}\}}{\operatorname{argmax}} \left(-1 + d(0.8) + e(0.2), 2 + d(0.1) + e(0.9) \right) \\ &= \underset{a \in \{a_{21}, a_{22}\}}{\operatorname{argmax}} \left(-1 + \frac{4d+e}{5}, 2 + \frac{d+9e}{10} \right) \end{aligned}$$

$$\text{now, } a^*(s_1) = a_{11} \text{ if } 5 + \frac{d+e}{2} > 10+e$$

$$\Rightarrow \frac{d-e}{2} > 5 \Rightarrow d-e > 10$$

$$\text{if } d-e < 10 \Rightarrow a^*(s_1) = a_{12}$$

$$\begin{aligned} a^*(s_2) &= a_{21} \text{ if } -1 + \frac{4d+e}{5} > 2 + \frac{d+9e}{10} \\ &\Rightarrow \frac{8d+2e}{10} > 3 + \frac{d+9e}{10} \\ &\Rightarrow 8d+2e > 30+d+9e \\ &\Rightarrow 7d > 30+7e \\ &\Rightarrow d-e > \frac{30}{7} \end{aligned}$$

$$a^*(s_2) = a_{22} \text{ if } d-e < \frac{30}{7}$$

$$\begin{aligned} (c) \gamma_1(s, d(s)) &= \begin{cases} p(A=a_{12}) \gamma_1(s_1, a_{12}) + p(A=a_{11}) \gamma_1(s_1, a_{11}) & ; s=s_1 \\ p(A=a_{21}) \gamma_1(s_2, a_{21}) + p(A=a_{22}) \gamma_1(s_2, a_{22}) & ; s=s_2 \end{cases} \\ &= \begin{cases} q[5] + (1-q)[10] & ; s=s_1 \\ 0[2] + 1[-1] & ; s=s_2 \end{cases} = \begin{cases} 10-5q & ; s=s_1 \\ -1 & ; s=s_2 \end{cases} \end{aligned}$$

$$P_1(j|S, d_1(S)) = \begin{cases} q \times [0.5] + (1-q)0 & ; j=S_1, S=S_1 \\ q \times [0.5] + (1-q)1 & ; j=S_2, S=S_1 \\ 1 \times 0.8 & ; j=S_1, S=S_2 \\ 1 \times 0.2 & ; j=S_2, S=S_2 \end{cases}$$

$$\Rightarrow P_1(j|S, d_1(S)) = \begin{cases} 0.5q & ; j=S_1, S=S_1 \\ 1 - 0.5q & ; j=S_2, S=S_1 \\ 0.8 & ; j=S_1, S=S_2 \\ 0.2 & ; j=S_2, S=S_2 \end{cases}$$

(d) Let $d: \mathcal{S} \rightarrow P(A)$ be a random policy s.t

$$d(S_1) = \begin{bmatrix} p \\ 1-p \end{bmatrix} \quad q_{11} \\ q_{12}$$

$$d(S_2) = \begin{bmatrix} q \\ 1-q \end{bmatrix} \quad a_{21} \\ a_{22}$$

then now to maximize:

$$TR(S_1) = p [5 + (5)(0.5) + (-5)(0.5)] \\ + (1-p)[10 + (5)(0) + (-5)(1)]$$

$$TR(S_2) = q [-1 + (5)(0.8) + (-5)(0.2)] \\ + (1-q)[2 + (5)(0.1) + (-5)(0.9)]$$

$$TR(S_1)(p) = p[5] + (1-p)(5) = 5$$

$$TR(S_2)(q) = q[-1 + 4 - 1] + (1-q)[2 - 4] \\ = q[2] + (1-q)(-2) \\ = 2q - 2 + 2q = 4q - 2$$

so, $d(S_1) = \begin{bmatrix} p \\ 1-p \end{bmatrix}$ does not depend on p
maximum value $q=1$

$$d(S_2) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

3. $\overbrace{\quad\quad\quad\quad\quad}$ digit number

$$T = \{1, 2, 3, 4, 5\} \quad N=5$$

at $\forall t \in T$

$$(a) X \sim \text{Unif}\{0, 1, 2, 3, \dots, 9\}$$

\downarrow_t
choose digit

(b) choose one of empty position

let positions be 1s, 10s, 100s, 1000s, 10000s

then Action = $\{1, 10, 100, 1000, 10000\} \rightarrow$ choose any 1
 $S = \text{state of current number, total places fixed}$

$$S_1 = (0, [00000]) \\ \hookrightarrow \text{No places filled}$$

$r_t(s_t, a) = x, x \leftarrow$ reward & we add to the number as total reward will optimise avg of number

where $A_{S_t} = \{10^{i-1} \mid i \text{th place } \& \text{ zero for digit for } s_t\}$

so if $s_t = (0, [00000])$

$$A_{S_1} = \{10^{1-1}, 10^{2-1}, \dots, 10^{5-1}\}$$

$$A_{S_1} = \{1, 10, 100, \dots, 10000\}$$

if $s_2 = (1, [00001])$

$$\text{then } A_{S_2} = \{10, 100, \dots, 10000\}$$

now, to formulate MDP:

$$(T, \mathcal{S}, A_S, S \in \mathcal{S}, r_t(s, a), P_t(j|s, a))$$

$$T = \{1, 2, \dots, 5\}$$

$$\mathcal{S} = \text{all states} = \{(\chi, [a_1, a_2, a_3, a_4, a_5]) \mid \chi = \sum_{i=1}^5 \alpha_i a_i, \alpha_i \in \{0, 1, \dots, 9\}, a_i \in \{0, 1\}^2\}$$

A_S for $s \in \mathcal{S}$ is

$$A_S = \{10^{i-1} \mid s = (\chi, [a_1, a_2, \dots, a_5]), i \text{ is s.t. } a_i = 1\}$$

$$s = (\chi, [a_1, \dots, a_5])$$

$$r_t(s, a) = \mathbb{E}[X_t] \times a$$

$$\mathbb{E}[X_t] = \frac{0+1+\dots+9}{10} = \frac{9 \times 10}{2 \times 10} = 4.5$$

$\Rightarrow r_t(s, a) = 4.5 \times a \rightarrow$ as total reward will optimise avg of final number formed

now, $P_t(j|s, a)$ now if $s = (\chi, [a_1, \dots, a_5])$

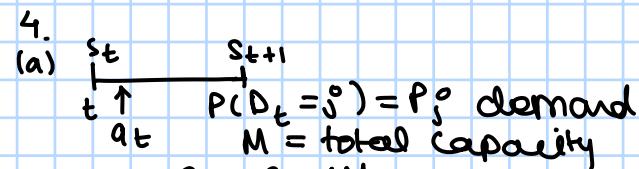
$$j = (y, [b_1, \dots, b_5])$$

now if action a is chosen then
 j will have $b_i = 1$ s.t. $a = 10^{i-1}$

$$\text{and } y = \chi + X_t a$$

$$\text{so, } P_t(j|s, a) = P_t(j|s, 10^{i-1}) = \begin{cases} \frac{1}{10}; j = (\chi + 10^{i-1}, [b_1, \dots, b_5]) \\ \downarrow a = 10^{i-1} \\ \text{s.t. } b_i = 1, b_j = 0, j \neq i \end{cases}$$

0; otherwise



$$a_t + s_t \leq M$$

$$r = \{1, \dots, N\}$$

$$j = \{0, 1, \dots, M\}$$

$$M = 3$$

$$N = 3 \text{ so } t=1, t=2, t=3 \rightarrow r_3(s_3) = 0 \text{ No action}$$

$$A_S = \{0, 1, \dots, M-1\}$$

$$s_{t+1} = \max \{0, s_t + a_t - D_t\}$$

$h(u)$ = holding cost
 $o(u)$ = ordering cost

$$o(u) = \begin{cases} 0 & ; u=0 \\ c(u) + k & ; u>0 \end{cases}$$

$\xrightarrow{\text{cost per } u}$ $\xrightarrow{\text{fixed cost}}$

$$o(u) = \begin{cases} 0 & ; u=0 \\ 4+2u & ; u\neq 0 \end{cases}$$

$f(j)$ = amount we get if we sell j units $f(u) = 8u$

$$r_t(s_t, a_t) = f(\min \underbrace{\{D_t, s_t + a_t\}}_{s_t + a_t - S_{t+1}}) - o(a_t) - h(s_t + a_t) \quad h(u) = u$$

$F(u) = \sum_{j=0}^{u-1} f(j) p(j) + f(u) [\sum_{j>u} p(j)] \rightarrow$ expected f value if we start with u units

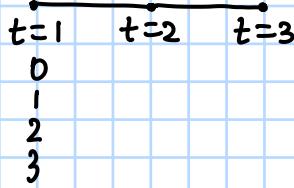
$$r_t(s_t, a_t) = F(s_t + a_t) - h(s_t + a_t) - o(a_t)$$

$$\begin{aligned} p_t(j|s, a) &= \begin{cases} 0 & ; j > s+a \\ p(D_t = s+a-j) & ; 0 < j \leq s+a \\ p(D_t > s+a) & ; j=0 \end{cases} \\ &= \begin{cases} 0 & ; j > s+a \\ p_{s+a-j} & ; 0 < j \leq s+a \\ \sum_i p_i^0 & ; j=0 \end{cases} \end{aligned}$$

$$P_0 = \frac{1}{4} \quad P_1 = \frac{1}{2} \quad P_2 = \frac{1}{4} \quad P_i^0 = 0 \text{ for } i \neq 0, 1, 2$$

$$\mathcal{S} = \{0, 1, 2, 3\}$$

Now,



$$\begin{aligned} A_0 &= \{0, 1, 2, 3\} \\ A_1 &= \{0, 1, 2\} \\ A_2 &= \{0, 1\} \\ A_3 &= \{0\} \end{aligned}$$

$$F(0) = 0$$

$$F(1) = P_0 \times f(0) + P_1 \times f(1) + P_2 \times f(1) = \left(\frac{1}{2} + \frac{1}{4}\right)(8) = 4 + 2 = 6$$

Start with 1

$$F(2) = P_0 \times f(0) + P_1 \times f(1) + P_2 \times f(2) = \frac{1}{2} \times 8 + \frac{1}{4} \times 2 = 8$$

$$F(3) = P_0 \times f(0) + P_1 \times f(1) + P_2 \times f(2) = 8$$

↑
Start with 3

now, $h(u) = u$ $o(u) = \begin{cases} 0 & ; u=0 \\ 4+2u & ; u\neq 0 \end{cases}$

$$r_t(s_t, a_t) = F(s_t + a_t) - h(a_t + s_t) - o(a_t)$$

$$r_3(s_3) = 0 \quad * s_3, a_3 \text{ possible as terminal case}$$

$$\gamma_1(0, 0) = F(0) - h(0) - u(0) = 0 - 0 - 0 = 0 = \gamma_2(0, 0)$$

$$\gamma_1(0, 1) = F(1) - h(1) - u(1) = 6 - 1 - 6 = -1 = \gamma_2(0, 1)$$

$$\gamma_1(0, 2) = F(2) - h(2) - u(2) = 8 - 2 - 4 = 2 = \gamma_2(0, 2)$$

$$\gamma_1(0, 3) = F(3) - h(3) - u(3) = 8 - 6 - 6 = -4 = \gamma_2(0, 3)$$

$$\begin{aligned} r_1(1,0) &= F(1) - h(1) - u(0) = 6 - 1 - 0 = 5 = r_2(1,0) \\ r_1(1,1) &= F(2) - h(2) - u(1) = 8 - 2 - 4 - 2 = 0 = r_2(1,1) \\ r_1(1,2) &= F(3) - h(3) - u(2) = 8 - 3 - 4 - 4 = -3 = r_2(1,2) \end{aligned}$$

$$\begin{aligned} r_1(2,0) &= F(2) - h(2) - u(0) = 8 - 2 - 0 = 6 = r_2(2,0) \\ r_1(2,1) &= F(3) - h(3) - u(1) = 8 - 3 - 4 - 2 = -1 = r_2(2,1) \\ r_1(3,0) &= F(3) - h(3) - u(0) = 8 - 3 = 5 = r_2(3,0) \end{aligned}$$

now, if at $t=2$ at state $s=0$

$$a_2^*(0) = \operatorname{argmax}_{a \in \{0, 1, 2, 3\}} (r_2(0, a) + \delta)$$

$$= 0 \quad \text{so, } a_2^*(0) = 0$$

$$a_2^*(1) = 0$$

$$a_2^*(2) = 0$$

$$a_2^*(3) = 0$$

so at $t=2$ we don't buy anything

now at $t=1$ $\delta_1 = 0$

$$a_1^*(0) = \operatorname{argmax} \left\{ r_1(0, a) + P(S_2=a) r_2(a, 0) + P(S_2=a-1) r_2(a-1, 0) + P(S_2=a-2) r_2(a-2, 0) \right\}$$

$$a_1^*(0): a=0 : r_1(0, 0) + 1 \times r_2(0, 0) = 0$$

$$a=1 : r_1(0, 1) + \frac{1}{4} r_2(1, 0) + \frac{3}{4} r_2(0, 0) = -1 + \frac{1}{4}(5) + 0 = \frac{1}{4}$$

$$\begin{aligned} a=2 : r_1(0, 2) + \frac{1}{4} r_2(2, 0) + \frac{1}{2} r_2(1, 0) + \frac{1}{4} r_2(0, 0) \\ = -2 + \frac{1}{4}(2) + \frac{1}{2}(5) + \frac{1}{4}(0) \\ = -2 + \frac{1}{4} = \frac{1}{4} \end{aligned}$$

$$\begin{aligned} a=3 : r_1(0, 3) + \frac{1}{4} r_2(3, 0) + \frac{1}{2} r_2(2, 0) + \frac{1}{4} r_2(1, 0) \\ = -8 + \frac{1}{4}(5) + \frac{1}{2}(6) + \frac{1}{4}(5) \\ = -8 + \frac{5}{2} + \frac{6}{2} = -8 + \frac{11}{2} \end{aligned}$$

$$a_1^*(0) = 2$$

$$a_1^*(1) : a=0 : r_1(1, 0) + \frac{1}{4} r_2(1, 0) + \frac{3}{4} r_2(0, 0) = 5 + \frac{1}{4}(5) + \frac{3}{4}(0) = 6.25$$

$$a=1 : r_1(1, 1) + \frac{1}{4} r_2(2, 0) + \frac{1}{2} r_2(1, 0) + \frac{1}{4} r_2(0, 0) = 0 + 2.5 = 2.5$$

$$a=2 : r_1(1, 2) + \frac{1}{4} r_2(3, 0) + \frac{1}{2} r_2(2, 0) + \frac{1}{4} r_2(1, 0) = -3 + 3 + 1.25 = 1.25$$

$$a_1^*(1) = 0$$

$$a_1^*(2) : a=0 : r_1(2, 0) + \frac{1}{4} r_2(2, 0) + \frac{1}{2} r_2(1, 0) + \frac{1}{4} r_2(0, 0)$$

$$= 6 + 1.5 + 2.5 + 0 = 10$$

$$a=1: \gamma_1(2,1) + \frac{1}{4} \gamma_2(3,0) + \frac{1}{2} \gamma_2(2,0) + \frac{1}{4} \gamma_2(1,0)$$

$$a_i^*(2) = -1 + 1.25 + 3 + 1.25$$

$$a_i^*(3) = 0 \rightarrow \text{trivial}$$

so at all states $a=0$, but if $s_i=0$ then $a_i=2$

(b) $b(u) = \text{backlogging cost}$, if applied at start

$$\text{i.e. } r(s_t, a_t) = f((s_t + a_t - D_t)^+) - o(a_t) - h(s_t + a_t) - b((D_t - s_t - a_t)^+)$$

for D_t
true

$$r(s_t, a_t) = \mathbb{E}[f(s_t + a_t - D_t)^+] - o(a_t) - h(s_t + a_t) - \mathbb{E}b(D_t - s_t - a_t)^+$$

$$\text{let } \mathbb{E}[f(u - D_t)^+] = F(u)$$

$$= \sum_{D_t=0}^{u-1} f(D_t) P(D_t) + \sum_{D_t=u}^{\infty} f(u) P(D_t)$$

$$\mathbb{E}[r(D_t - u)^+] = R(u)$$

$$= \sum_{D_t=0}^u r(0) P(D_t) + \sum_{D_t=u+1}^{\infty} r(D_t - u) P(D_t)$$

$$\text{then } r_t(s_t, a_t) = F(s_t + a_t) - o(a_t) - h(s_t + a_t) - r(s_t + a_t)$$

$$P(D_t=0) = y_1 \quad P(D_t=1) = y_2 \quad P(D_t=2) = y_3$$

$$\text{and } u_t^*(s_t) = \max_{a_t \in A_{s_t}} (r(s_t, a_t) + \sum_{j \in S} p(j|s_t, a_t) u_{t+1}^*(j))$$

Reward-to-go function

$$a_t^*(s_t) = \operatorname{argmax}_{a_t \in A_{s_t}} (r(s_t, a_t) + \sum_{j \in S} p(j|s_t, a_t) u_{t+1}^*(j))$$

Optimal quantity to order at t, given s_t

$$\text{now } R(u) = \sum_{D_t=0}^u r(0) P(D_t) + \sum_{D_t=u+1}^{\infty} r(D_t - u) P(D_t)$$

for $u=0$

$$R(0) = \sum_{D_t=1}^2 r(D_t) P(D_t)$$

$$= \frac{1}{2} \times 3 + \frac{1}{4} \times 3 \times 2$$

$$R(0) = 3$$

$$R(1) = \frac{1}{4} \times 3 \times 1 = \frac{3}{4}$$

$$R(2) = 0$$

$R(3) = 0$, every other function same as (a)

(() code on website

5. S, A s are countable

$N = \text{Horizon length} (\ N < \infty)$

(a) π be HR policy i.e. $\pi = (d_1, d_2, \dots, d_{N-1})$

s.t.

$$d_i : H_i \rightarrow P(A)$$

↑
takes in history till i and outputs probability
of action

$$H_i = S_i \times A_1 \times \dots \times S_{i-1} \times A_{i-1} \times S_i$$

$$V_N^\pi(S) = \mathbb{E}_S^\pi \left[\sum_{t=1}^{N-1} r_t(x_t, a_t) + \gamma_N(x_N) \right]$$

start at S

$U_t^\pi : H_t \rightarrow \mathbb{R}$ be reward-to-go function

$$U_t^\pi(h_t) = \mathbb{E}_{h_t}^\pi \left[\sum_{n=t}^{N-1} r_n(x_n, a_n) + \gamma_N(x_N) \right]$$

$$= r_t(s_t, d_t(h_t)) + \mathbb{E}_{h_t}^\pi \left[\sum_{n=t+1}^{N-1} r_n(x_n, a_n) + \gamma_N(x_N) \right]$$

$$U_t^\pi(h_t) = r_t(s_t, d_t(h_t))$$

$$+ \mathbb{E}_{h_t}^\pi \left[\mathbb{E}_{h_{t+1}}^\pi \left(\sum_{n=t+1}^{N-1} r_n(x_n, a_n) + \gamma_N(x_N) \right) \right]$$

by tower property

$$= r_t(s_t, d_t(h_t)) + \mathbb{E}_{h_t}^\pi [U_{t+1}^\pi(h_{t+1}, d_{t+1}(h_t), x_{t+1})]$$

only thing
that is random

now algorithm: set $t=N$, $U_N^\pi(h_N) = \gamma_N(s_N)$ $\forall h_N \in H_N$

then put $t=N-1$:

$$U_{N-1}^\pi(h_{N-1}) = \gamma_{N-1}(s_{N-1}, d_{N-1}(h_{N-1})) \quad \forall h_{N-1} \in H_{N-1}$$

$$+ \mathbb{E}_{h_{N-1}}^\pi [U_N^\pi(h_{N-1}, d_{N-1}(h_{N-1}), x_N)]$$

iterate $t=N-2, N-3, \dots, 1$

$$\text{then } U_1^\pi(S) = V_N^\pi(S) \quad \forall S \in S$$

$$(b) \text{ Now } U_t(h_t) = \sup_{a \in A_{S_t}} \left\{ r_t(s_t, a) + \sum_{j \in S} p_t(j|s_t, a) U_{t+1}(h_t, a, j) \right\}$$

$U_N(h_N) = r_N(s_N)$ is the bellman equation where

U_t^* be the solution to bellman equation, it is also
the optimal strategy to follow

now, $U_t^*(h_t)$ only depends on h_t through s_t as

$$\text{for } t=N: U_N^*(h_N) = \gamma_N(s_N) \quad (\text{true for } t=N)$$

if true for $t \in \{N-1, \dots, 1\}$ then for n :

$$U_n^*(h_n) = \sup_{a \in A_{S_n}} \left\{ r_n(s_n, a) + \sum_{j \in S} p_n(j|s_n, a) U_{n+1}^*(j) \right\}$$

→ true for n

now, as solution to bellman equation depends on s_t
and not a_t

$$U_n^*(s_t) = \sup_{a \in A_{s_t}} \left\{ r_t(s_t, a) + \sum_{s' \in S} p(s'|s_t, a) U_{t+1}^*(s') \right\}$$

by defn of supremum, $\exists a \in A_{s_t}$ s.t for given ϵ
 $\epsilon + U_n^\pi(s_t) \geq U_n^*(s_t)$

then action a depends on s_t for given t
and so similarly for every time same

$$\epsilon + U_1^\pi(s_1) \geq U_1^*(s_1)$$

$$\pi = (d_1, d_2, \dots, d_{N-1})$$

$$\text{s.t } d_i : S_i \rightarrow A_{S_i} \text{ thus } d_i(s_i) = a_i \text{ s.t } \epsilon\text{-optimal}$$

(c) MR optimal policy $\pi = (d_1, \dots, d_{N-1})$ given

$$\text{s.t } d_i : S_i \rightarrow P(A_{S_i}) \quad \text{this is MD from (b)}$$

$$\text{optimal policy } \pi \text{ s.t } U_t^*(s_t) = \sup_{\pi \in \text{MR}} U_t^\pi(s_t)$$

and the supremum is for given $\pi \in \text{MR}$

then π is given optimal MD policy, now

$$\text{MD policy } \pi' = (d'_1, \dots, d'_{N-1})$$

$$d'_i : S_i \rightarrow A_{S_i}$$

optimal policy should satisfy bellman
equation, we can find

$U_t^\pi(s_t)$ using (a) algorithm and so

$$\text{let } a_t^*(s_t) = \underset{a \in A_{s_t}}{\operatorname{argmax}} \left(r_t(s_t, a) + \sum_{s' \in S} p(s'|s_t, a) U_{t+1}^*(s') \right)$$

as same as
 $U_{t+1}^*(s')$
and computed
by our algorithm

then $a_t^*(s_t)$ is optimal action to take at
time t given s_t

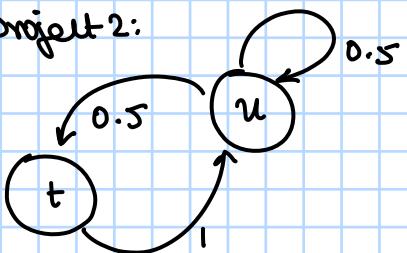
$$\text{or } d'_i : S \rightarrow A_S$$

$$d'_i(s) = a_t^*(s)$$

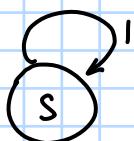
$$d'_i(s) = \underset{a \in A_s}{\operatorname{argmax}} \left(r_t(s, a) + \sum_{s' \in S} p(s'|s, a) U_{t+1}^*(s') \right)$$

$\pi' = (d'_1, d'_2, \dots, d'_{N-1})$ is optimal MD policy

6. project-2:



project-1:



now, $t=1$ 2 cases $(S, t), (S, u)$
 $t=2$ 2 cases $(S, t), (S, u)$
 $t=3$ 2 cases $(S, t), (S, u)$

so, total states = $\{(S, t), (S, u)\} = \{$

$$\text{actions: } A_{(S, t)} = \{1, 2\}$$

$$A_{(S, u)} = \{1, 2\} \leftarrow \text{what project to pick}$$

$$\gamma((S, t), 1) = 1 \quad \gamma((S, t), 2) = 0$$

$$\gamma((S, u), 1) = 1 \quad \gamma((S, u), 2) = 2$$

$$\gamma_3(S) = 0 \quad \forall S \in \{$$

$$\begin{aligned} \text{now, } p((S, t) | (S, t), 1) &= 1 \\ p((S, t) | (S, u), 1) &= 0 \\ p((S, u) | (S, t), 1) &= 0 \\ p((S, u) | (S, u), 1) &= 1 \\ p((S, t) | (S, t), 2) &= 0 \\ p((S, u) | (S, t), 2) &= 1 \\ p((S, t) | (S, u), 2) &= 1/2 \\ p((S, u) | (S, u), 2) &= 1/2 \end{aligned}$$

now, for $N=3$:

$$u_3^*((S, t)) = 0$$

$$u_3^*((S, u)) = 0$$

$N=2$:

$$\begin{aligned} u_2^*(S, t) &= \max \left\{ 1 + \sum_{i=1}^2 () (0), 0 + \sum_{i=1}^2 () (0) \right\} \\ &= \max \{1, 0\} \end{aligned}$$

$$u_2^*(S, t) = 1 \quad \text{so } a_2^*(S, t) = 1$$

$$u_2^*(S, u) = \max_{i=2} \left\{ 1 + \sum_{j=1}^{i-1} () (0), 2 + \sum_{j=1}^{i-1} () (0) \right\}$$

$$u_2^*(S, u) = 2, \quad a_2^*(S, u) = 2$$

$$\text{now, } u_1^*(S, t) = \max \{ 1 + 1 \times 1, 0 + 1 \times 2 \}$$

$$= 2 \quad a_1^*(S, t) = \{1, 2\} \quad \text{does not matter}$$

$$u_1^*(s, u) = \max \{ 1 + 1 \times 2, 2 + \gamma_2 \times 2 + \gamma_2 \times 1 \}$$

$$= \max_{=3.5} \{ 3, 3.5 \}$$

$$q_1^*(s, u) = 2$$

so, State action at $t=1$ action at $t=2$

(s, u)	2	2
(s, t)	1 or 2	1

Tutorial-2:

$$1. \quad \mathbb{V}_y^{\pi}(S) = \mathbb{E}_S^{\pi} \left[\mathbb{E}_y \left[\sum_{t=1}^{\gamma} \gamma(X_t, Y_t) \right] \right]$$

$\gamma \sim \text{Hypergeometric}$

$$P(\gamma = n) = (n-1) (1-\lambda)^2 \lambda^{n-2} \quad n=2, 3, \dots$$

$$\text{To prove: } \mathbb{V}_y^{\pi}(S) = \mathbb{V}_{\lambda}^{\pi}(S) + (1-\lambda) \frac{\partial \mathbb{V}_{\lambda}^{\pi}(S)}{\partial \lambda}$$

$$\text{proof: } \mathbb{V}_{\lambda}^{\pi}(S) = \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{t=1}^N \lambda^{t-1} \gamma(X_t, Y_t) \right]$$

$$\frac{\partial \mathbb{V}_{\lambda}^{\pi}(S)}{\partial \lambda} = \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{t=1}^N (t-1) \lambda^{t-2} \gamma(X_t, Y_t) \right]$$

$$\Rightarrow (1-\lambda) \frac{\partial \mathbb{V}_{\lambda}^{\pi}(S)}{\partial \lambda} = \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{t=1}^N (1-\lambda)(t-1) \lambda^{t-2} \gamma(X_t, Y_t) \right]$$

$$\Rightarrow \mathbb{V}_y^{\pi}(S) + (1-\lambda) \frac{\partial \mathbb{V}_{\lambda}^{\pi}(S)}{\partial \lambda} = \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{t=1}^N \lambda^{t-1} \gamma(X_t, Y_t) + \sum_{t=1}^N (1-\lambda)(t-1) \lambda^{t-2} \gamma(X_t, Y_t) \right]$$

$$= \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{t=1}^N (\lambda^{t-1} + (1-\lambda)(t-1) \lambda^{t-2}) \gamma(X_t, Y_t) \right] \quad \textcircled{1}$$

$$\begin{aligned} \text{now, } \mathbb{V}_y^{\pi}(S) &= \mathbb{E}_S^{\pi} \left[\mathbb{E}_y \left[\sum_{t=1}^{\gamma} \gamma(X_t, Y_t) \right] \right] \\ &= \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{n=1}^N P(\gamma = n) \left(\sum_{t=1}^n \gamma(X_t, Y_t) \right) \right] \\ &= \lim_{N \rightarrow \infty} \mathbb{E}_S^{\pi} \left[\sum_{n=1}^N (n-1)(1-\lambda)^2 \lambda^{n-2} \left(\sum_{t=1}^n \gamma(X_t, Y_t) \right) \right] \\ &= \mathbb{E}_S^{\pi} \left[\sum_{n=1}^{\infty} (n-1)(1-\lambda)^2 \lambda^{n-2} \left(\sum_{t=1}^n \gamma(X_t, Y_t) \right) \right] \\ &= \mathbb{E}_S^{\pi} \left[\sum_{t=1}^{\infty} \left(\underbrace{\sum_{n=t}^{\infty} (n-1)(1-\lambda)^2 \lambda^{n-2}}_{\textcircled{2}} \right) \gamma(X_t, Y_t) \right] \quad \textcircled{3} \\ &\quad (\because \text{telescopic sum}) \end{aligned}$$

$$\text{Now for all A.P of type: } ab + (a+d)b\gamma + (a+2d)b\gamma^2 + \dots = S$$

given $r < 1 \quad rS = ab\gamma + (a+d)b\gamma^2 + \dots$

$$\begin{aligned} S - rS &= ab + db\gamma + db\gamma^2 + db\gamma^3 + \dots \\ (1-r)S &= ab + db\gamma \left(\frac{1}{1-r} \right) \end{aligned}$$

$$S = \frac{ab}{1-r} + \frac{db\gamma}{(1-r)^2}$$

$$\begin{aligned} \text{so, } \textcircled{2} \text{ is s.t } \left. \begin{aligned} a &= t^{-1} \\ b &= (1-\lambda)^2 \lambda^{t-2} \\ d &= 1 \\ r &= \lambda \end{aligned} \right\} \quad S = (t-1)(1-\lambda) \lambda^{t-2} + \lambda^{t-1} \end{aligned}$$

$$\begin{aligned}
 80, \quad \vartheta_{\nu}^{\pi}(s) &= \mathbb{E}_s^{\pi} \left[\sum_{t=1}^{\infty} ((t-1)(1-\lambda)\lambda t^{-2} + \lambda t^{-1}) r(x_t, y_t) \right] \\
 &= \lim_{N \rightarrow \infty} \mathbb{E}_s^{\pi} \left[\sum_{t=1}^N \lambda t^{-1} + (t-1)(1-\lambda)\lambda t^{-2} \right] r(x_t, y_t)
 \end{aligned}$$

$$\vartheta_{\nu}^{\pi}(s) = \vartheta_{\lambda}^{\pi}(s) + (1-\lambda) \frac{\partial \vartheta_{\lambda}^{\pi}(s)}{\partial \lambda} \quad \text{from } ①$$

2. $\mathcal{S} = \{1, 2, \dots\}$ $A_S = \{a\}$ $P(s+1|s, a) = 1$, $r(s, a) = 0 \quad \forall s, a$

$$(a) \quad \vartheta_1 = (0, 0, \dots)$$

$$\vartheta_2 = (1, \frac{1}{\lambda}, \frac{1}{\lambda^2}, \dots)$$

$$d(s) = a \Rightarrow r_d = (r(s, d(s))) = (0, 0, 0, \dots)$$

$$(P_d)_{s,j} = P(j|s, d(s)) = \begin{cases} 1 & j = s+1 \\ 0 & \text{otherwise} \end{cases}$$

$$P_d = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & \dots \end{bmatrix}$$

$$\begin{aligned}
 L\vartheta &= \sup_a \{ r_d + \lambda P_d \vartheta \} \\
 &\quad \vartheta_1 = (0, 0, \dots) = r_d \\
 L\vartheta_1 &= \sup_a \{ \vartheta_1 + \lambda \begin{bmatrix} 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & \dots \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ \vdots \end{bmatrix} \} \\
 &= (0, 0, 0, \dots)
 \end{aligned}$$

$$L\vartheta_2 = \lambda \begin{bmatrix} 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix} \begin{bmatrix} 1 \\ \frac{1}{\lambda} \\ \frac{1}{\lambda^2} \\ \frac{1}{\lambda^3} \\ \vdots \end{bmatrix} = \lambda \begin{bmatrix} 1 \\ \frac{1}{\lambda} \\ \frac{1}{\lambda^2} \\ \frac{1}{\lambda^3} \\ \vdots \end{bmatrix} = \vartheta_2$$

$$80, \quad L\vartheta_1 = \vartheta_1, \quad L\vartheta_2 = \vartheta_2$$

$$(b) \quad \text{let } \mathcal{V} = \{ \vartheta \mid \|\vartheta\| = \sup_{s \in \{1, 2, \dots\}} \vartheta(s), < \infty \}$$

then as we define $L: \mathcal{V} \rightarrow \mathcal{V}$
as d is only one

$$\begin{aligned}
 L\vartheta &= r_d + \lambda P_d \vartheta \\
 &= \lambda P_d \vartheta
 \end{aligned}$$

$$\begin{aligned}
 L(\vartheta^1, \vartheta^2, \dots) &= \lambda(\vartheta^2, \vartheta^3, \dots) \\
 \text{now } L(\vartheta^1, \vartheta^2, \dots) - L(\vartheta^1, \vartheta^2, \dots) &= \lambda((\vartheta^2, \vartheta^3, \dots) - (\vartheta^2, \vartheta^3, \dots))
 \end{aligned}$$

$$= \lambda \sup_{s \in \{2, 3, \dots\}} \|\vartheta^s - \vartheta^s\|$$

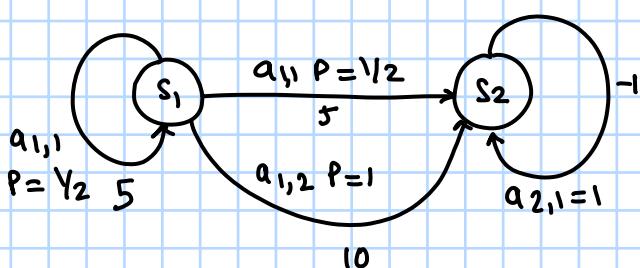
but if $\mathbf{v}, \mathbf{u} \in \mathcal{V}$ $\|\mathbf{u} - \mathbf{v}\| = \|\mathbf{u}' - \mathbf{v}'\|$ and $\mathbf{u}', \mathbf{v}' \in \mathcal{V}$

$$\|u^l - v^l\| > \|u^s - v^s\| \forall s \in \{2, \dots\}$$

true

$\|Lu - Lv\| > \lambda \|v - u\|$ so L is not a contraction map
and we can't apply banach fixed point theorem
(if v is banach, if v is not banach then also same)

3.



We want to maximise reward to go

$$V(S_2) = r(S_2, a_{2,1}) + \lambda P(S_2 | S_2, a_{2,1}) V(S_2)$$

↑
Reward to go for S_2 at time t

↑
Reward to go at $t+1$

$$\Rightarrow (1-\lambda) V(S_2) = -1$$

$$\Rightarrow V(S_2) = \frac{-1}{1-\lambda}$$

now for S_1 : $V(S_1) = \underbrace{\max_{a \in A_S} r(S_1, a) + \lambda (P(S_2 | S_1, a) V(S_2) + P(S_1 | S_1, a) V(S_1))}_{\text{Reward to go at } t}$

$$\begin{aligned} &= \max \left\{ r(S_1, a_{1,1}) + \lambda \left(P(S_2 | S_1, a_{1,1}) V(S_2) + P(S_1 | S_1, a_{1,1}) V(S_1) \right), \right. \\ &\quad \left. r(S_1, a_{1,2}) + \lambda P(S_2 | S_1, a_{1,2}) V(S_2) \right\} \\ &= \max \left\{ \underbrace{5 + \lambda \left(\frac{1}{2} V(S_2) + \frac{1}{2} V(S_1) \right)}_{a_{1,1}}, \underbrace{10 + \lambda V(S_2)}_{a_{1,2}} \right\} \end{aligned}$$

$$V(S_2) = \frac{-1}{1-\lambda}$$

$$V(S_1) = \max \left\{ \underbrace{s + \frac{\lambda}{2} \left(V(S_1) - \frac{1}{1-\lambda} \right)}_{a_{1,1}}, \underbrace{10 - \frac{\lambda}{1-\lambda}}_{a_{1,2}} \right\}$$

now $a^*(S_1) = a_{1,1}$ if

$$\textcircled{1} \quad s + \frac{\lambda}{2} \left(V(S_1) - \frac{1}{1-\lambda} \right) > 10 - \frac{\lambda}{1-\lambda}$$

$$\textcircled{2} \quad V(S_1) = s + \frac{\lambda}{2} \left(V(S_1) - \frac{1}{1-\lambda} \right) \Rightarrow \left(1 - \frac{\lambda}{2} \right) V(S_1) = s - \frac{\lambda}{2(1-\lambda)}$$

$a^*(S_1) = a_{1,2}$ if

$$\Rightarrow V(S_1) = \frac{s - \frac{\lambda}{2(1-\lambda)}}{1 - \frac{\lambda}{2}}$$

$$\textcircled{1} \quad s + \frac{\lambda}{2} \left(V(S_1) - \frac{1}{1-\lambda} \right) < 10 - \frac{\lambda}{1-\lambda}$$

$$\textcircled{2} \quad V(S_1) = 10 - \frac{\lambda}{1-\lambda}$$

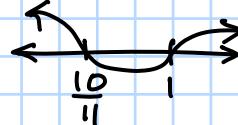
so, for $a^*(s_1) = a_{1,1}$ we get:

$$\begin{aligned} \frac{s+\lambda}{2} \left(\frac{s-\lambda}{\lambda(1-\lambda)} - \frac{1}{1-\lambda} \right) &> 10 - \frac{\lambda}{1-\lambda} \\ \frac{\lambda}{2} \left(\frac{10}{2-\lambda} - \frac{2\lambda}{\lambda(1-\lambda)(2-\lambda)} - \frac{1}{1-\lambda} \right) &> s - \frac{\lambda}{1-\lambda} \\ \Rightarrow \frac{\lambda}{2} \left(\frac{10(1-\lambda)-(2-\lambda)}{(2-\lambda)(1-\lambda)} - \frac{\lambda}{(1-\lambda)(2-\lambda)} \right) &> \frac{s(1-\lambda)-\lambda}{1-\lambda} \\ \Rightarrow \frac{\lambda}{2} \left(\frac{10-10\lambda-2+\lambda-\lambda}{2-\lambda} \right) &> s - 5\lambda - \lambda \\ \Rightarrow \lambda(5-5\lambda-1) &> s - 6\lambda \\ \Rightarrow \lambda \left(\frac{4-5\lambda}{2-\lambda} \right) &> s - 6\lambda \\ \Rightarrow 4\lambda - 5\lambda^2 &> 10 - 12\lambda - 5\lambda + 6\lambda^2 \\ \Rightarrow 11\lambda^2 - 21\lambda + 10 &< 0 \end{aligned}$$

$$\lambda = \frac{21 \pm \sqrt{(21)^2 - 4(11)(10)}}{2 \times 11}$$

$$= \frac{21 \pm 1}{22}$$

$$\text{so, } (\lambda-1)(\lambda-\frac{10}{11}) < 0$$



$\therefore \lambda \in \left(\frac{10}{11}, 1 \right)$ for $a^*(s_1) = a_{1,1}$

$$\text{now, } s + \frac{\lambda}{2} \left(10 - \frac{\lambda}{1-\lambda} - \frac{1}{1-\lambda} \right) < 10 - \frac{\lambda}{1-\lambda}$$

$$\Rightarrow \frac{\lambda}{2} \left(\frac{10-10\lambda-\lambda-1}{1-\lambda} \right) < \frac{s-5\lambda-\lambda}{1-\lambda}$$

$$\Rightarrow \frac{\lambda}{2} (9-11\lambda) < s - 6\lambda$$

$$\Rightarrow 9\lambda - 11\lambda^2 < 10 - 12\lambda$$

$$\Rightarrow 0 < 11\lambda^2 - 21\lambda + 10$$

$$\Rightarrow (\lambda-1)(\lambda-\frac{10}{11}) > 0$$

$$\lambda \in [0, \frac{10}{11}]$$

for $a^*(s_1) = a_{1,2}$

4. $L: \mathcal{T} \rightarrow \mathcal{T}$

To prove: L is a contraction map

Proof: let $u, v \in \mathcal{T}$
and let $Lu = \sup_{d \in DMD} \{r_d + \lambda p_d u\}$ by defn

Goal: $L\vartheta(s) \geq L\psi(s)$ for some $s \in S$

$$L\vartheta(s) = \sup_{a \in A_s} \left\{ \gamma(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta(j) \right\}$$

by defn of sup, $\forall \varepsilon > 0, \exists a^* \in S$ s.t.

$$\gamma(s, a^*) + \lambda \sum_{j \in S} P(j|s, a^*) \vartheta(j) + \varepsilon \geq L\vartheta(s)$$

$$\text{and } L\psi(s) \leq \gamma(s, a^*) + \lambda \sum_{j \in S} P(j|s, a^*) \psi(j)$$

$$L\vartheta(s) - L\psi(s) \leq \varepsilon + \lambda \sum_{j \in S} P(j|s, a^*) (\vartheta(j) - \psi(j))$$

$$\Rightarrow |L\vartheta(s) - L\psi(s)| \leq \varepsilon + \lambda \left\| \sum_{j \in S} P(j|s, a^*) \right\| \|(\vartheta - \psi)\|$$

$$\Rightarrow \|L\vartheta - L\psi\| \leq \varepsilon + \lambda \|\vartheta - \psi\| \quad (\because \text{true} \forall s \in S)$$

$$\Rightarrow \|L\vartheta - L\psi\| \leq \lambda \|\vartheta - \psi\|$$

$$\text{so, } \forall u, v \in V, \|Lu - Lv\| \leq \lambda \|u - v\|$$

so L is a contraction map

$$5. \text{ To prove: } \sup_{d \in D^{MD}} \{r_d + \lambda p_d \vartheta\} = \sup_{d \in D^{MR}} \{r_d + \lambda p_d \vartheta\}$$

Proof:

$$\text{Let } F\vartheta(s) = \sup_{\substack{a \in A_s \\ j \in S}} \left\{ \gamma(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta(j) \right\}$$

we define this function

$$\tilde{F}\vartheta(s) = \sup_{P(A_s)} \sum_{a \in A_s} \left(P(a) \gamma(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta(j) \right)$$

now to maximize weighted average, we have to assign all weights to one action that is maximum
so from this

$$\sup_{a \in A_s} \left\{ \gamma(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta(j) \right\} = \sup_{P(A_s)} \sum_{a \in A_s} P(a) \gamma(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta(j)$$

$$\Rightarrow F(\vartheta(s)) = \tilde{F}(\vartheta(s)) \quad \forall s \in S$$

$$\Rightarrow \sup_{d \in D^{MD}} \{r_d + \lambda p_d \vartheta\} = \sup_{d \in D^{MR}} \{r_d + \lambda p_d \vartheta\}$$

6. To prove: If $v \in V$ s.t. $Lv = v$ then

$$v = \sup_{\pi \in \Pi^{SR}} v^\pi$$

Proof: If $v \geq Lv$ then (comp wise)

now as $Lv \geq r_d + \lambda p_d v$ by defn (comp wise)

$$\Rightarrow v \geq r_d + \lambda p_d v$$

and $v_\lambda^\pi = r_d + \lambda P_d v_\lambda^\pi$ for $\pi \in \Pi^{SR}$ (\because done in class)

$$(I - \lambda P_d) v \geq r_d$$

as $(I - \lambda P_d)^{-1}$ is positive (\because done in class)

$$\Rightarrow v \geq (I - \lambda P_d)^{-1} r_d$$

$$v_\lambda^\pi \text{ for } \pi = (d, d, \dots)$$

$$\Rightarrow v \geq v_\lambda^\pi \text{ for } \pi \in \Pi^{SR}$$

as this is true $\forall d$, we get

$$\forall \pi \in \Pi^{SR} \Rightarrow v \geq v_\lambda^\pi$$

$$\Rightarrow v \geq v_\lambda^* = \sup_{\pi \in \Pi^{SR}} v_\lambda^\pi \quad \text{--- ①}$$

now, if $v \leq L v$ then from defn of $\sup_{\forall \epsilon > 0, \exists d \in D^{SD} \text{ s.t.}}$

$$\begin{aligned} v &\leq r_d + \lambda P_d v + \epsilon \cdot e \\ \Rightarrow (I - \lambda P_d) v &\leq r_d + \epsilon \cdot e \\ \Rightarrow v &\leq (I - \lambda P_d)^{-1} r_d + (I - \lambda P_d)^{-1} \epsilon \cdot e \\ &= v_\lambda^\pi \text{ for } \pi \in \Pi^{SD} \quad \frac{\epsilon \cdot e}{1-\lambda} \end{aligned}$$

$$\Rightarrow v \leq v_\lambda^\pi + \frac{\epsilon \cdot e}{1-\lambda}$$

taking sup

$$\Rightarrow v \leq \sup_{\pi \in \Pi^{SD}} v_\lambda^\pi + \frac{\epsilon \cdot e}{1-\lambda} \leq \sup_{\pi \in \Pi^{SR}} \underbrace{v_\lambda^\pi}_{\text{as } \Pi^{SD} \subseteq \Pi^{SR}} + \frac{\epsilon \cdot e}{1-\lambda}$$

as $\Pi^{SD} \subseteq \Pi^{SR}$

$$\Rightarrow \sup_{\pi \in \Pi^{SD}} v_\lambda^\pi \leq \sup_{\pi \in \Pi^{SR}} v_\lambda^\pi$$

$$\Rightarrow v \leq \sup_{\pi \in \Pi^{SR}} v_\lambda^\pi \text{ as } \epsilon \rightarrow 0 \quad \text{--- ②}$$

now if $v = L v$ then $v \leq L v \Rightarrow v \leq v_\lambda^*$ from ②
 $v = L v$ then $v \geq L v \Rightarrow v \geq v_\lambda^*$ from ①

$$\begin{aligned} \text{so as } v &\leq v_\lambda^* \text{ & } v \geq v_\lambda^* \\ \Rightarrow v &= v_\lambda^* \\ \Rightarrow v &= \sup_{\pi \in \Pi^{SR}} v_\lambda^\pi \end{aligned}$$

7. given $\exists \pi \in \Pi^{HR}$ s.t.

$$v_\lambda^\pi = v_\lambda^* = \sup_{\pi \in \Pi^{HR}} v_\lambda^\pi$$

To prove: $\exists \tilde{\pi} \in \Pi^{SD}$ s.t. $v_\lambda^\pi = v_\lambda^{\tilde{\pi}}$
proof: let $\pi = (d_1, d_2, \dots) \in \Pi^{HR}$, then we know
 $\exists \pi' \in (d'_1, d'_2, \dots) \in \Pi^{MR}$ s.t.

$$d_t': \mathcal{S} \rightarrow P(A)$$

$$P\pi'(x_n=j, y_n=a | x_1=s) = P\pi(x_n=j, y_n=a | x_1=s) \quad (\text{theorem done})$$

i.e. $V_{\lambda}^{\pi'}(s) = V_{\lambda}^{\pi}(s)$ in class (d_2', d_3', \dots)

$$V_{\lambda}^{\pi'}(s) = \sup_{P(A)} \sum_{a \in A_s} \sum_{j \in \mathcal{S}} \gamma(s, a) q_{d_t'(s)}(a) + \lambda \sum_{j \in \mathcal{S}} \left(\sum_{a \in A_j} P(j | s, a) q_{d_t'(s)}(a) \right) V_{\lambda}^{\pi}(j)$$

$$\text{where } q_{d_t'(s)}(j) = P(y_t=a | x_t=j, x_1=s)$$

as sup over $P(A_s)$, $\exists a \in A_s$, s.t. $\{P(A_s=a)=1\}$ \Rightarrow argmax of all $P(A_s)$

as average maximum, mean all weights given to

so let this a be a_s^* then

$$V_{\lambda}^{\pi'}(s) = \gamma(s, a_s^*) + \lambda \sum_{j \in \mathcal{S}} P(j | s, a_s^*) V_{\lambda}^{\pi}(j) \quad (d_2'(j), \dots)$$

let $d: \mathcal{S} \rightarrow A$ s.t.

$d(s) = a_s^*$ as above can be done HS
now,

$$V_{\lambda}(j) \leq V_{\lambda}(j) \uparrow^{(d_1'(j), \dots)} = V_{\lambda}^{\pi}(j)$$

$$\Rightarrow V_{\lambda}(j) \leq r(j, a_j^*) + \lambda \sum_{K \in \mathcal{S}} P(K | j, a_j^*) V_{\lambda}(K) \quad (\text{as optimal policy})$$

$$\Rightarrow V_{\lambda}(j) \leq r(j, a_j^*) + \lambda \sum_{K \in \mathcal{S}} P(K | j, a_j^*)$$

$$(r(k, a_k^*) + \lambda \sum_{K_2 \in \mathcal{S}} P(K_2 | k, a_k^*))$$

$$(r(k_2, a_{k_2}^*) + \lambda \sum_{K_3 \in \mathcal{S}} \dots)$$

$$\Rightarrow V_{\lambda}(j) \leq \sum_{K=0}^{\infty} \lambda P_d^K r_d(j) = V_{\lambda}(d, d, \dots)(j)$$

$$\Rightarrow V_{\lambda}(j) \leq V_{\lambda}(d, d, \dots)(j)$$

$$\text{now, } V_{\lambda}^{\pi'}(s) \leq r(s, d(s)) + \lambda \sum_{j \in \mathcal{S}} P(j | s, d(s)) V_{\lambda}(d, \dots)(j)$$

$$V_{\lambda}^{\pi'}(s) \leq V_{\lambda}(s)$$

$$\text{but as } V_{\lambda}^{\pi'}(s) = \sup_{\pi \in \Pi_{HR}} V_{\lambda}^{\pi} \leq V_{\lambda}(s)$$

$$\Rightarrow V_{\lambda}^{\pi'}(s) = V_{\lambda}(s)$$

$$\text{so, let } \tilde{\pi} = (d, d, \dots) \quad V_{\lambda}^{\pi'}(s) = V_{\lambda}^{\tilde{\pi}}(s) = V_{\lambda}^{\pi}(s) \quad \forall s \in \mathcal{S}$$

so, V_{λ}^{π} is optimal

8. $\theta \in V$ $d \in D^{MD}$ is θ -improving if $d \in \operatorname{argmax}_{d \in D^{MD}} \{ r_d + \lambda P_d \theta \}$

(a) $\vartheta_\lambda^{(d_V, d_V, \dots)} < v$ given $d_V \in \operatorname{argmax}_{d \in D^{MD}} \{ r_d + \lambda P_d \theta \}$

Let $v = (1000)$ only one state s
 $r(s, a) = (0)$ only one action a

then $d_V(s) = a$ only one $d \in D^{MD}$

$$\vartheta_\lambda^{\pi} = (I - \lambda P_d)^{-1} \vartheta_d = (0) \Rightarrow \vartheta_\lambda^{\pi} = (0)$$

$$\vartheta_\lambda^{\pi} < v$$

(b) Let S, A only one state and one action and let $r_d = (1-\lambda)$

then

$$\text{for } \theta = (1)$$

we get

$$v = (1) \leq (1-\lambda) + \lambda(1) = (1)$$

and

$$\vartheta_\lambda^{(d_V, d_V, \dots)} = (I - \lambda P_{d_V})^{-1} \vartheta_{d_V} = \frac{1}{1-\lambda} (1-\lambda) = (1)$$

$$\text{so, } r_{d_V} + \lambda P_{d_V} v \geq v \text{ and } \vartheta_\lambda^{(d_V, d_V, \dots)} = v$$

(c) To prove: $\vartheta_\lambda^{(d_V, \dots)}(s'') < v(s'')$ for some $s'' \in S$ if $r_{d_V}(s') + \lambda P_{d_V} v(s') \geq v(s')$ for some $s' \in S$

Proof: $\exists s' \in S \quad r_{d_V}(s') + \lambda P_{d_V} v(s') \geq v(s')$

as $d_V \in \operatorname{argmax}_{d \in D^{MD}} \{ r_d + \lambda P_d \theta \}$

$$r_{d_V}(s') + \lambda P_{d_V} v(s') \leq r_{d_V}(s) + \lambda P_{d_V} v(s) \quad \forall d \in D^{MD}$$

if $\exists s'' \in S \text{ s.t. } r_{d_V}(s'') + \lambda P_{d_V} v(s'') < v(s'')$

then, $r_{d_V}(s'') + \lambda P_{d_V} v(s'') < v(s'') \quad \forall d \in D^{MD}$

$$\Rightarrow (I - \lambda P_d)^{-1} r_{d_V}(s'') < v(s'') \quad \forall d \in D^{MD}$$

$$\Rightarrow (I - \lambda P_{d_V})^{-1} r_{d_V}(s'') < v(s'')$$

$$\Rightarrow \vartheta_\lambda^{(d_V, d_V, \dots)}(s'') < v(s'')$$

if $\forall s \in S \quad r_{d_V}(s) + \lambda P_{d_V} v(s) \geq v(s)$

$$\Rightarrow r_{d_V} + \lambda P_{d_V} v \geq v \quad (\text{component wise})$$

$$\Rightarrow (I - \lambda P_{d_V})^{-1} r_{d_V} \geq v$$

$\Rightarrow \vartheta_\lambda^{(d_V, d_V, \dots)} \geq v$, so as in statement $\exists s'$, we get
 for all other $s'' \neq s'$ we get above, if
 all s' , then not possible

9. Finite state and actions $\{S_1, S_2, \dots, S_N\}$ $N \geq 1$ $d \in D^{MD}$ let $P_d^L + P_d^U = P_d$

$$P_d^L = (I - \lambda P_d^U)$$

$$P_d^U = \lambda P_d V$$

$$T: V \longrightarrow V$$

$$TV = \max_{d \in DMD} \{ Q_d^\top r_d + Q_d^\top R_d v \}$$

$$\text{wlog } TV(s) > TU(s)$$

$$\text{let } d^* \in \arg\max_{d \in DMD} \{ Q_d^\top r_d + Q_d^\top R_d v \}$$

$$TV(s) = Q_{d^*}^\top r_{d^*}(s) + Q_{d^*}^\top R_{d^*} v(s)$$

$$\text{and } TU(s) \geq Q_{d^*}^\top r_{d^*}(s) + Q_{d^*}^\top R_{d^*} u(s)$$

$$\Rightarrow TV(s) - TU(s) \leq Q_{d^*}^\top r_{d^*}(s) + Q_{d^*}^\top R_{d^*} v(s)$$

$$- Q_{d^*}^\top r_{d^*}(s) - Q_{d^*}^\top R_{d^*} u(s)$$

$$\Rightarrow TV(s) - TU(s) \leq Q_{d^*}^\top R_{d^*} (v(s) - u(s))$$

$$\Rightarrow |TV(s) - TU(s)| \leq \|Q_{d^*}^\top R_{d^*}\| |v(s) - u(s)| \quad \forall s \in S$$

$$\Rightarrow \|TV - TU\| \leq \|Q_{d^*}^\top R_{d^*}\| \|v - u\|$$

now seen in class $\|Q_{d^*}^\top R_{d^*}\| \leq 1$, so T is a contraction map

10. on website, nothing on monotonicity as r not ≥ 0
or ≤ 0

