

30th Sept:

standard LP: $\min \sum \alpha(s) \vartheta(s)$ where $\vartheta(s) \geq g_i(s, a) + \lambda \sum_{j \in S} p(j|s, a) \vartheta(j)$
 $\frac{\vartheta}{\alpha}$ is the primal, $\alpha(s) > 0 \forall s$, $\sum \alpha(s) = 1 \forall s, a$
 from above LP we get optimal value

dual: $\max \sum_{s,a} \gamma(s, a) \alpha(s, a)$ s.t. $\sum_{a \in S^0} \gamma(j, a) - \lambda \sum_{s,a} p(j|s, a) \alpha(s, a) = \alpha(j) \forall j$

this LP gives optimal policy $(d^*)^\infty$, \rightarrow feasible region

Note: If $\sum_{s,a} c(s, a) \alpha(s, a) \leq \beta$ \forall for constraint MDP

now, $\lim_{T \rightarrow \infty} \sum_{k=1}^T \sum_{t=1}^T p^{d^{\infty}}(x_t=s, y_t=a | x_1=k)$
 $= \mathbb{E} [N(x_t=s, y_t=a \text{ in } t=1, \dots, T \text{ steps})]$ If this was taken certain
it will converge to stationary
state

$\sum_{k=1}^{\infty} \alpha(k) \sum_{t=1}^T \lambda^{t-1} p^{d^{\infty}}(x_t=s, y_t=a | x_1=k)$ is discounted occupancy measure
 $= \chi_d(s, a) \leftarrow \infty$, $\chi(s, a)$ is just discounted occupancy measure

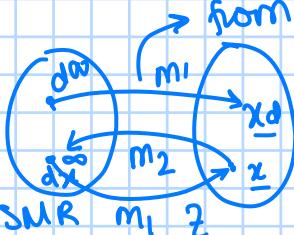
$$\begin{aligned} \sum_{s,a} \chi_d(s, a) \gamma(s, a) &= \sum_{s,a} \alpha(k) \sum_{t=1}^{\infty} \lambda^{t-1} \sum_{s,a} p^{d^{\infty}}(x_t=s, y_t=a | x_1=k) \gamma(s, a) \\ &= \sum_{s,a} \alpha(k) \sum_{t=1}^{\infty} \lambda^{t-1} [\mathbb{E}[g(x_t, y_t | x_1=k)]] \quad \text{different notations} \\ &= \sum_s \vartheta^{d^{\infty}}(s) \alpha(s) = \vartheta^{d^{\infty}}(\alpha) = J^{d^{\infty}}_{\alpha, \lambda} = J(d^{\infty}, \alpha) \end{aligned}$$

from below theorem, (a) $\Rightarrow \vartheta^{d^{\infty}}(\alpha) = \sum_{s,a} g(s, a) \chi_d(s, a)$ as $\underline{\chi}_d \in \mathcal{Z}$ from below theorem

$$(b) \Rightarrow \sum_{s,a} g(s, a) \chi(s, a) = \vartheta^{d^{\infty}}(\alpha)$$

as $\underline{\chi} = \underline{\chi}_d$ we have

$$\begin{aligned} \sum_{s,a} \chi(s, a) \gamma(s, a) &= \sum_{s,a} \chi_{d,\underline{\chi}}(s, a) \gamma(s, a) \\ &= \vartheta^{d^{\infty}}(\underline{\chi}) \end{aligned}$$



from (b) given any $d \rightarrow \underline{\chi}_d$ given any $\underline{\chi} \rightarrow d^*$
 M_1, M_2 are mappings so one-one, onto

∞ , max of feasible region gives us optimal policy
Note: $r(s, a)$ is not special here, we can use $c(s, a)$ to do same for constraint MDP

Karlin's result: $\chi^*(s, a) = 0 \nmid a \neq a_s^*, s \neq s^*$

$\chi^*(s^*, a^*) > 0, \chi^*(s^*, a^*_{s^*}) > 0$ when constraint

Theorem: (a) If $d^{\infty} \in \text{SMR}$, then define

$$\chi_d(s, a) \triangleq \sum_{k \in S} \alpha(k) \sum_{t=1}^{\infty} \lambda^{t-1} p^{d^{\infty}}(x_t=s, y_t=a | x_1=k)$$

then $\underline{\chi}_d \in \mathcal{Z}$

$$(b) \text{ If } \underline{\chi} \in \mathcal{Z} \text{ define } d_{\underline{\chi}}(s, a) = \frac{\chi(s, a)}{\sum_{a'} \chi(s, a')} \quad \forall s, a$$

now $d_{\underline{\chi}} \in \text{SMR}$

Let

$$\underline{\chi}_{d_{\underline{\chi}}} \text{ be s.t. } \chi_{d_{\underline{\chi}}}(s, a) \triangleq \sum_{k \in S} \alpha(k) \sum_{t=1}^{\infty} \lambda^{t-1} p^{d_{\underline{\chi}}}(x_t=s, y_t=a | x_1=k)$$

then $\underline{\chi}_{d_{\underline{\chi}}} = \underline{\chi}$

Proof: (a) given d^∞ , defined \underline{x}_d , we have to show $\underline{x}_d \in \mathcal{Z}$

$$\begin{aligned}
 \text{now } & \lambda \sum_{S,a} p(j|s,a) x_d(s,a) && \text{(part LHS from dual LP)} \\
 & = \lambda \sum_{S,a} p(j|s,a) \sum_{k=1}^{\infty} \lambda^{t-1} p^{d^\infty}(x_t=s, y_t=a | x_1=k) \\
 & = \lambda \sum_k \sum_{t=1}^{\infty} \lambda^{t-1} \sum_{S,a} p(j|s,a) p^{d^\infty}(x_t=s, y_t=a | x_t=k) \\
 & = \sum_k \sum_{t=1}^{\infty} \lambda^t \sum_a p^{d^\infty}(x_{t+1}=j, y_{t+1}=a | x_1=k) \\
 & = \sum_k \sum_{n=2}^{\infty} \lambda^{n-1} \sum_a p^{d^\infty}(x_n=j, y_n=a | x_1=k) \\
 & = \sum_k \lambda^n \sum_a p^{d^\infty}(x_n=j, y_n=a | x_1=k) \\
 & \quad - \sum_k \lambda^n \sum_a p^{d^\infty}(x_1=j, y_1=a | x_1=k) \\
 & = \sum_a x_d(j,a) - \alpha(j) \sum_k \delta(j=k) \quad \text{will be 1}
 \end{aligned}$$

so, we get: $\sum_a x_d(j,a) - \alpha(j) = \lambda \sum_{S,a} p(j|s,a) x(s,a) + j$

this is same as $\underline{x} \in \mathcal{Z}$ from dual LP

One: prove (b) of theorem \rightarrow done down

Theorem: (i) The two LPs have optimal solution

(ii) say \underline{x}^* , \underline{v}^* are some optimal solution then

(a) $\underline{v}^*(s)$ value vector of MDP

(b) $d^*(s,a) = \frac{x^*(s,a)}{\sum_a x^*(s,a)}$ then $(d^*)^\infty$ is an optimal policy for MDP

assuming $|x(s,a)| \leq M$

Proof: Feasibility will be true as for any $d \in \text{SMR}$ from (a) of above theorem we get $\underline{x}_d \in \mathcal{Z}$

$$\begin{aligned}
 \text{now as } & \sum_{S,a} x_d(s,a) v(s,a) = \sum_{K \in S} \alpha(K) v^{d^\infty}(s) \leq \frac{M}{1-\lambda} \\
 & \text{as } -\frac{M}{1-\lambda} \leq v^{d^\infty}(s) \leq \frac{M}{1-\lambda}
 \end{aligned}$$

$$\text{so, } -\frac{M}{1-\lambda} \leq \sum_{S,a} x_d(s,a) v(s,a) \leq \frac{M}{1-\lambda}$$

so its feasible, bounded and so max will exist and so min min

(i) of theorem follows from this as feasible and bounded

(ii) we get \underline{x}^* , \underline{v}^* are optimal solution

$$\sum a(s) v^*(s) = \sum x^*(s,a) \alpha(s,a) \quad (\because \text{duality})$$

$$\Rightarrow \sum \alpha(s) \underline{v}^*(s) = \sum \pi^*(s, a) \underline{v}(s, a)$$

from construction, we can use \underline{v} for \underline{v}^*
 $= \sum \alpha(s) \underline{v} d_{\underline{x}^*}^{\infty}(s)$ from previous theorem
(b) part

now we have to show $\underline{v}^*(s) = \underline{v} d_{\underline{x}^*}^{\infty}(s) \forall s$, then we are done
using below lemma we have

$$\underline{v}^*(s) \geq \underline{v} d_{\underline{x}^*}^{\infty}(s) \quad (\because d_{\underline{x}^*}^{\infty} \text{ is also a policy})$$

now as $\sum \alpha(s) \underline{v}^*(s) = \sum \alpha(s) \underline{v} d_{\underline{x}^*}^{\infty}(s)$
and $\alpha(s) > 0$

$$\Rightarrow \sum \alpha(s) \underline{v}^*(s) \geq \sum \alpha(s) \underline{v} d_{\underline{x}^*}^{\infty}(s)$$

now, as $\sum \alpha(s) \underline{v}^*(s) = \sum \alpha(s) \underline{v} d_{\underline{x}^*}^{\infty}(s)$
if $\exists s.t. \underline{v}(s) > \underline{v} d_{\underline{x}^*}^{\infty}(s)$
then $\sum \alpha(s) \underline{v}^*(s) > \sum \alpha(s) \underline{v} d_{\underline{x}^*}^{\infty}(s)$
this is a contradiction

$$\therefore \underline{v}^*(s) = \underline{v} d_{\underline{x}^*}^{\infty}(s) \forall s$$

$d_{\underline{x}^*}^{\infty}$ is optimal policy for MDP

Lemma: $\underline{v}^*(s) \geq \underline{v} d^{\infty}(s) \forall s, \forall d^{\infty}$ where $\underline{v}^*(\cdot)$ is solution to primal

Proof: as \underline{v}^* is solution to primal it is in \mathcal{S} i.e

$$\underline{v}^*(s) \geq r(s, a) + \lambda \sum_{j \in \mathcal{S}} p(j|s, a) \underline{v}^*(j) \quad \forall s, a$$

now for any $d^{\infty} \in \text{SMR}$ we have

$$\sum_a q_d(s, a) = 1 \quad (\because \text{probability distribution})$$

$$\Rightarrow \underline{v}^*(s) = \sum \underline{v}^*(s) q_d(s, a)$$

$$\geq \sum_a (r(s, a) + \lambda \sum_{j \in \mathcal{S}} p(j|s, a) \underline{v}^*(j)) q_d(s, a)$$

$$= \sum_a [q_d(s, a) r(s, a)] + \sum_a [\lambda \sum_{j \in \mathcal{S}} p(j|s, a) q_d(s, a) \underline{v}^*(j)]$$

$$= r_d(s) + \lambda \sum_{j \in \mathcal{S}} p_d(j|s) \underline{v}^*(j) \quad \forall s$$

$$\Rightarrow \underline{v}^* \geq r_d + \lambda p_d \underline{v}^*$$

$$\text{as } \underline{v} d^{\infty} = r_d + \lambda p_d \underline{v} d^{\infty}$$

$$\Rightarrow \underline{v}^* \geq \underline{v} d^{\infty} - \lambda p_d \underline{v} d^{\infty} + \lambda p_d \underline{v}^*$$

$$\Rightarrow \underline{v}^*(I - \lambda p_d) \geq \underline{v} d^{\infty}(I - \lambda p_d)$$

$$\Rightarrow \underline{v}^* \geq \underline{v} d^{\infty} \quad (\because \text{seen before, in Positive maps})$$

Ex: Show (b) of theorem

AM: $\pi(s, a) \in \mathbb{Z}$ let

$$u(s) = \sum_a \pi(s, a) \quad \text{true } u(s) > 0 \quad (\text{from dual constraints})$$

and if d_x is valid $\Rightarrow d_x^{\infty}$ is valid SMR

now to show $x = x_d$

$$x(j) = u(j) - \sum_{s, a} \lambda p(j|s, a) \pi(s, a)$$

$$= u(j) - \sum_{s, a} \lambda p(j|s, a) \pi(s, a) \frac{u(s)}{\sum \pi(s, a)}$$

$$= u(j) - \sum_{s,a} \lambda P(j|s,a) q_d(s,a) u(s)$$

$$= u(j) - \sum_s \lambda P(j|s) u(s)$$

i.e. $\underbrace{(\alpha(s_1), \dots, \alpha(s_n))}_{\alpha^T} = (\underbrace{u(s_1), \dots, u(s_n)}_{u^T}) - \lambda(u(s_1), \dots, u(s_n)) \times \begin{bmatrix} P_{d,x}(s_1|s_1), \dots \\ \vdots \\ P_{d,x}(s_n|s_n), \dots \end{bmatrix}$

$$\Rightarrow \alpha^T = u^T - \lambda u^T P_{d,x}$$

$$= u^T(I - \lambda P_{d,x})$$

$$\Rightarrow \alpha^T (I - \lambda P_{d,x})^{-1} = u^T$$

$$\Rightarrow \alpha^T \left(\sum_{n=1}^{\infty} (\lambda P_{d,x})^n \right)^{-1} = u^T$$

so, $u(s) = \sum_j \alpha(j) \sum_{n=1}^{\infty} \lambda^n \sum_{a \in As} p_{d,x}^{(n)}(x_n=s, y_n=a | x_1=j)$

$$= \sum_{a \in As} x_{d,x}(s, a)$$

$$\Rightarrow \sum_{a \in As} x(s, a) = \sum_{a \in As} x_{d,x}(s, a) \quad \text{--- ①}$$

now, $x_{d,x}(s, a) = \sum_j \alpha(j) \sum_{n=1}^{\infty} \lambda^n p_{d,x}^{(n)}(x_n=s | x_1=j) q_{d,x}(s)(a)$

$$= \sum_j \alpha(j) \sum_{n=1}^{\infty} \lambda^n p_{d,x}^{(n)}(x_n=s | x_1=j) \frac{x(s, a)}{\sum_{a' \in As} x(s, a')} \quad \text{--- ②}$$

and $\sum_a x_{d,x}(s, a) = \sum_j \alpha(j) \sum_{n=1}^{\infty} \lambda^n p_{d,x}^{(n)}(x_n=s | x_1=j) \quad \text{--- ③}$

so, ② becomes: $x_{d,x}(s, a) = \left(\sum_{a_1} x_{d,x}(s, a_1) \right) \left(\frac{x(s, a)}{\sum_{a_2} x(s, a_2)} \right) \quad (\because ③)$

$$= \left(\cancel{\sum_{a_1} x_{d,x}(s, a_1)} \right) (x(s, a)) \quad (\because ①)$$

$$\cancel{\left(\sum_{a_2} x_{d,x}(s, a_2) \right)}$$

so, $x_{d,x}(s, a) = \cancel{x(s, a)} + s, a$

$$\Rightarrow \underline{x_{d,x}} = \underline{x}$$

3rd Oct:

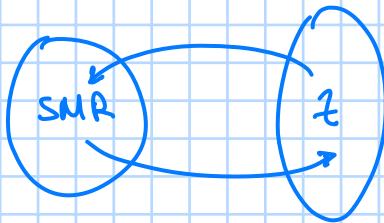
- ① $\min_{\underline{\alpha}} \sum_s \alpha(s) \underline{V}(s)$, $\underline{V}(s) > \pi(s, a) + \lambda \sum_j P(j|s, a) \underline{V}(j) \quad \forall s, a$
- ② $\max_{\underline{\pi}} \sum_{s, a} \pi(s, a) \pi(s, a) \quad \sum_a \pi(s', a) = \lambda \sum_{s, a} P(s'|s, a) \pi(s, a) + \pi(s') \quad \forall s'$

we also saw $\forall z \in \mathbb{Z}$, $\exists d_x(s, a) = \frac{\pi(s, a)}{\sum_{a'} \pi(s, a')}$
 s.t $d_x \in \text{SMR}$

and similarly $\forall d \in \text{SMR}, \exists \underline{x} \in \mathbb{Z}$

$$x_d(s, a) = \sum_k \alpha(k) \sum_{t=1}^{\infty} \lambda^{t-1} P^{\infty}(X_t=s, Y_t=a | X_1=k)$$

so, given any vector $\underline{z} \in \mathbb{Z}$ $\exists f \in \text{SMR}$



$\exists f: \text{SMR} \rightarrow \mathbb{R}$
 f is one-one and onto
 and so, $\exists f^{-1}$

we have seen, $\forall z, \exists d_x^{\infty}$ and $\forall d^{\infty}, \exists \underline{x} \in \mathbb{Z}$

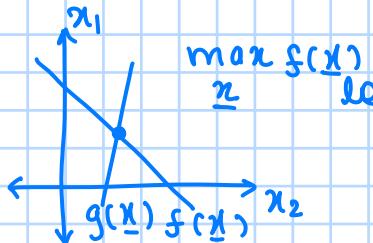
also using d_x^{∞} , if we define $\underline{x}d_x$ we get

Note: $\Omega_{\pi}(x) = \sum_{s, a} \pi(s, a) \pi(s, a)$

$$\text{then, } \underline{V}^{\pi} \text{ is s.t } \underline{V}^{\pi}(x) = \Omega_{\pi}(x)$$

$$\sum_s \alpha(s) \underline{V}^{\pi}(s)$$

Note: If we have constraint MDP, $\sum x(s, a) c(s, a) \leq \beta$ will represent constrained function



$$\max_{\underline{x}} f(\underline{x}) \text{ s.t. } x_1 + x_2 = \beta$$

$$\text{let } g(\underline{x}) \leq \beta$$

constraint

$$\begin{aligned} & \left(\sup_{\pi} \sum_{s, t} \lambda^{t-1} E_{\pi}^{\pi}[R_t] \right. \\ & \left. \text{s.t. } \sum_{s, t} \lambda^{t-1} E_{\pi}^{\pi}[c_t] \leq \beta \forall \pi \right) \\ & \text{wait to find } \underline{\pi} \end{aligned}$$

Note: we also saw if \underline{V}^* , \underline{x}^* are solution to ① and ②, then we saw given $\underline{x}^* \rightarrow d_x^{\infty} \in \text{SMR}$

and $\underline{V}^* = \underline{V}^{\pi^*}$ for any value of π i.e. it does not depend on π

we only need

$$\begin{aligned} \alpha(s) &> 0 \quad \forall s \\ \sum_s \alpha(s) &= 1 \end{aligned}$$

we get this as there is potential to optimise all of them in every direction (as $\alpha(s) > 0$) and so does not matter on α

any d^{∞} , we can show $\underline{V}^* \geq \underline{V}^{\infty}$

$$\Rightarrow \underline{V}^* \geq \sup_{\pi} \underline{V}^{\pi} = \underline{V}^*$$

$\Rightarrow \underline{V}^* \geq \underline{V}_x^*$ Point wise

and $\underline{V}^* \geq \underline{V}^{d\pi^\infty}$ as for every $d^\infty \in \text{SMR}$
 $\Rightarrow \underline{V}^*(s) \geq \underline{V}^{d\pi^\infty}(s) \quad \forall s$

$\underline{V}^{d\pi^\infty}$ satisfies ① with equality

and \underline{V}^* ,

$\underline{V}^{d\pi^\infty} \in \text{feasible primal LP}$
 and as \underline{V}^* is optimal

$$\text{so, } \sum_s \alpha(s) \underline{V}^*(s) \leq \sum_s \alpha(s) \underline{V}^{d\pi^\infty}(s)$$

as minimal value as \underline{V}^* is optimal

if $\exists s$ s.t. $\underline{V}^*(s) > \underline{V}^{d\pi^\infty}(s)$
 then

$$\Rightarrow \sum_s \alpha(s) \underline{V}^*(s) < \sum_s \alpha(s) \underline{V}^{d\pi^\infty}(s)$$

$\Rightarrow 1 < 1$ this is a contradiction

$$\text{so, } \underline{V}^*(s) = \underline{V}^{d\pi^\infty}(s) \quad \forall s$$

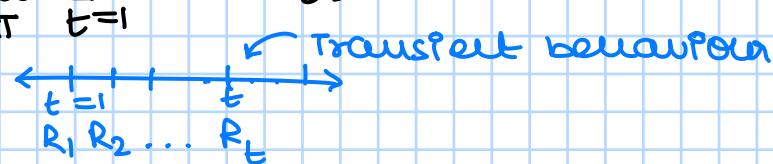
so, finally we get $\underline{V}^* = \underline{V}^{d\pi^*} = \underline{V}_\lambda^*$

$$\text{as } \underline{V}^* \geq \underline{V}_\lambda^* \Rightarrow \underline{V}^{d\pi^\infty} \geq \underline{V}_\lambda^* \\ \Rightarrow \underline{V}^{d\pi^\infty} = \underline{V}_\lambda^* \text{ as } d\pi^\infty \in \text{SMR}$$

so, $d\pi^*$ is optimal policy

(also $\underline{V}^* = \underline{V}_\lambda^*$ and so \underline{V}^* does not depend on λ)
 \rightarrow so by \underline{V}_λ^* we can get optimal policy

Now, we saw discounted cost MDPs where we want to maximize
 $\sup_{\pi} \sum_{t=1}^{\infty} \lambda^{t-1} \mathbb{E}^{\pi}[R_t]$



Eg: for a queuing system we are seeing something for 12 hours
 we keep working at rejection cost every 5 min
 here $\lambda=1$ to model an MDP as reputation on future
 interval is important ($\lambda=1$) \leftarrow lot of importance to study state

Average cost MDPs:

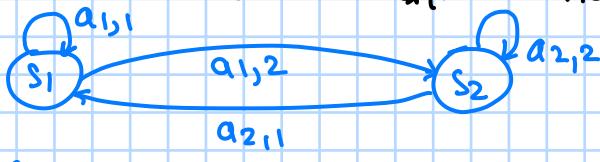
we now put $\lambda=1$, and try to maximize $\lim_{T \rightarrow \infty} \frac{\sum_{t=1}^T \mathbb{E}^{\pi}[R_t]}{T}$

also, $\exists \bar{\lambda} \in (0,1)$ st. $\bar{\lambda}$ is some big value

$$d_\lambda^* = d_{\bar{\lambda}}^* \text{ average cost } \forall \lambda > \bar{\lambda}$$

$$\text{i.e. } (1-\lambda)\underline{V}_{*,\lambda} \rightarrow \underline{V}_{*,\text{avg}}$$

Eg:



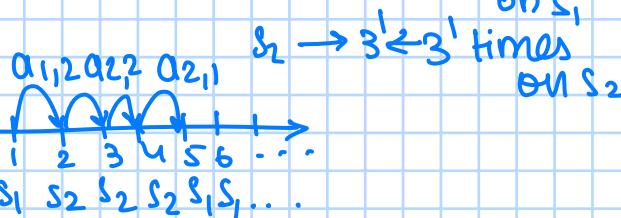
$$P(S_1 | S_1, a_{11}) = 1 \\ P(S_1 | S_1, a_{12}) = 1$$

$$P(S_2 | S_2, a_{22}) = 1 \\ P(S_2 | S_2, a_{21}) = 0$$

$$\sigma(S_1, a) = 2 \quad \forall a \in A_{S_1} \\ \sigma(S_2, a) = -2 \quad \forall a \in A_{S_2}$$

now if we start with s_1 , always

$$\pi \in HD, s_1 \rightarrow 3^0 \leftarrow 3^0 \text{ times}$$



$$s_1 \rightarrow 3^2 \text{ time}$$

$$s_1 \rightarrow 3^{2k}$$

$s_2 \rightarrow 3^{2k+2}$, from at end of $2k$ switches, $2k+1$ switches

$$N_k = \sum_{n=0}^{2k} (3)^n$$

$$= \frac{3^{2k+1}-1}{3-1}$$

$$N_k + 3^{2k+1}$$

$$N_k = \sum_{n=0}^{2k+1} (3)^n$$

now, let $\frac{v_{N_k}}{N_k} = \text{avg over } N_k \text{ steps}$

$$\begin{aligned} \frac{v_{N_k}}{N_k} &= \frac{\sum_{n=0}^{k-1} 3^{2n} (2) - \sum_{n=0}^{k-1} 3^{2n+1}}{N_k} \\ &= \frac{2}{N_k} \left(\sum_{n=0}^{k-1} 3^{2n} - \sum_{n=0}^{k-1} 3^{2n+1} \right) \\ &= \frac{4}{3^{2k+1}-1} \left(\frac{(3^2)^{k+1}-1}{3^2-1} - \frac{(3^2)^{k-1}}{3^2-1} \right) \end{aligned}$$

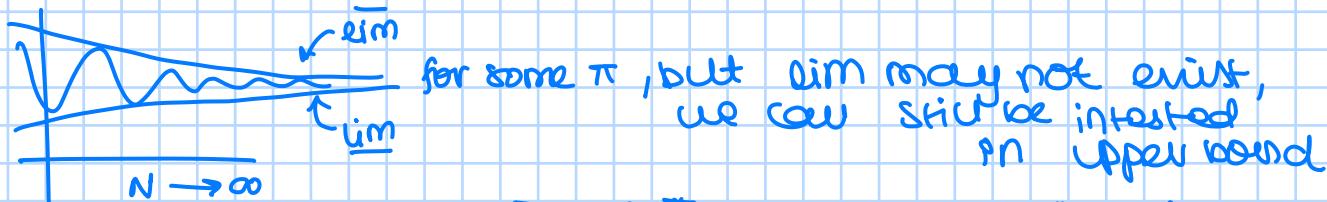
$$\text{for } k \rightarrow \infty \text{ we get } \lim_{k \rightarrow \infty} \frac{v_{N_k}}{N_k} = \frac{1}{2}$$

similarly for $N_k + 3^{2k+1}$ part

$$\lim_{k \rightarrow \infty} \frac{v_{N_k+3^{2k+1}}}{N_k+3^{2k+1}} = -\frac{1}{2}$$

$$\text{so, } \overline{\lim}_{N \rightarrow \infty} \frac{v_N}{N} = \frac{1}{2} \quad \underline{\lim}_{N \rightarrow \infty} \frac{v_N}{N} = -\frac{1}{2} \Rightarrow \text{not convergent}$$

Note: From above example we see that not all policy will converge to avg



i.e. $\inf_{\pi} \overline{\lim}_{N \rightarrow \infty} \frac{v_N}{N}^{\pi}$ or $\sup_{\pi} \underline{\lim}_{N \rightarrow \infty} \frac{v_N}{N}^{\pi} \rightarrow$ depending on use case of MDP

now, $x \in S$ $P(S'|S, d(S))$ some markov chain by fixing d

Defn: (irreducible) $\exists n$ s.t. $P^n(S'|S) > 0$ if $S \rightarrow S'$, $\forall S, S'$
return in finite time $x \rightarrow x \rightarrow \dots$

Defn: If $P_n(\tau_{X<\infty}) = 1$, $\tau_x \rightarrow$ return time

$$\tau_x = \inf \{n \geq 1; X_n = x\}$$

then x is recurrent

Note: If irreducible, and $\exists x \in S$ s.t. x is recurrent $\Rightarrow \forall s \in S$ s is recurrent

Defn: If $E_x[\tau_x] < \infty$ then x is called positive recurrent

Note: If irreducible, and $\exists x$ s.t. x is positive recurrent \Rightarrow all positive recurrent

now, if we have a markov chain which is positive recurrent

$$\Leftrightarrow \exists S, D \text{ s.t. } \pi = \pi P_d$$

$\pi(x) = \frac{1}{E_x[\tau_x]}$
 π is dist, $\pi(x) \rightarrow$ prob of being in x , $x \in S$

$$\pi P_d(s) = \sum_s \pi(s) P_d(s'|s)$$

$$P(X_2=s|X_1 \sim \pi)$$

now, $\frac{\sum_{n=1}^N \mathbb{I}\{X_n=s\}}{N} \rightarrow \pi(s)$ almost surely
as $N \rightarrow \infty$

and also $\frac{E\left[\sum_{n=1}^N \mathbb{I}\{X_n=s\}\right]}{N} \rightarrow \pi(s)$
as $N \rightarrow \infty$

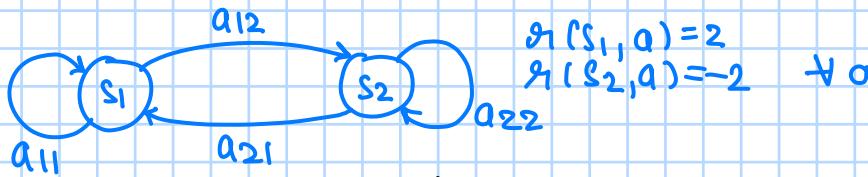
i.e. $\frac{\sum_{n=1}^N P_M(X_n=s)}{N} \xrightarrow{\text{Expected value}} \pi(s) \quad \forall M$
as $N \rightarrow \infty$

we want to look at $\frac{\sum_{n=1}^N E[\pi(X_n)]}{N}$
 $= \frac{\sum_S \sum_{n=1}^N P_M(X_n=s) \pi(s)}{N}$
 $\rightarrow \sum_S \pi(s) \pi(s)$

use M : $\sum_{S_i} \pi(S_i) P(X_n=S_i | X_1=S_i) = P_M(X_n=S_i)$
initial distribution

7th Oct:

Eg:



$$g_1(s_1, a) = 2$$

$$g_1(s_2, a) = -2 \quad \forall a$$

$$P(s_j | s_i, a_{ij}) = 1$$

$$P(s_i | s_i, a_{ii}) = 0$$

We are interested in $\lim_{N \rightarrow \infty} \frac{\sum_{n=1}^N E[g(x_n, y_n)]}{N}$ if it exists, in above

example, we can make a policy which is history dependent

$$x_0 = s_1, \quad a \xrightarrow{\begin{matrix} 1 \\ 3 \\ 3^2 \dots \end{matrix}} \begin{matrix} s_1 \\ s_2 \\ s_1 \end{matrix}$$

This is like pseudo stability

$s_1 \quad s_2 \quad s_1 \quad \dots$ were we saw

$$\overline{\lim}_{n=1}^N E[g(x_n, y_n)] = 2$$

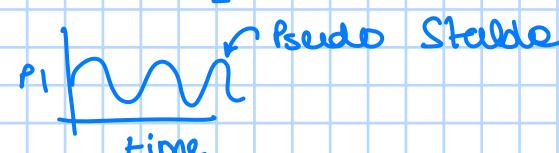
$$\underline{\lim}_{n=1}^N \frac{E[g(x_n, y_n)]}{N} = -2$$

Eg: Let P_1 = proportion of type 1 bacteria

If P_1 is small then $\rightarrow 0$



but if P_1, P_2 large both $\frac{P_1}{P_2} = 1$, then as time $\rightarrow \infty$



any policy $\pi \in \Pi^{HR}$, we want to find

$$\bar{V}^\pi(s) = \lim_{N \rightarrow \infty} \frac{\sum_{n=1}^N E_s^\pi [g(x_n, y_n)]}{N}$$

this may not exist

and then $\sup_\pi \bar{V}^\pi(s)$

We know, $\pi \in \Pi^{HR}$ and $s = x_0$, $\exists \pi' \in \Pi^{MR}$ such that all marginal distribution are the same

Only at time epoch n

$$\pi \leftrightarrow (x_n, y_n) \stackrel{\text{def}}{=} \pi' \leftrightarrow (x'_n, y'_n) \quad \forall n$$

$$\therefore V_N^\pi(s) = \sum_{n=1}^N E_s^\pi [g(x_n, y_n)] = V_N^{\pi'}(s) \quad \forall N$$

$$\Rightarrow \overline{\lim}_{N \rightarrow \infty} \frac{V_N^\pi(s)}{N} = \overline{\lim}_{N \rightarrow \infty} \frac{V_N^{\pi'}(s)}{N} \quad \text{and similarly} \quad \lim_{N \rightarrow \infty} \frac{V_N^\pi(s)}{N} = \lim_{N \rightarrow \infty} \frac{V_N^{\pi'}(s)}{N}$$

If both $\overline{\lim}$ and $\underline{\lim}$ same then above limit will exist

$$\text{Defn: } g_+^\pi(s) = \lim_{N \rightarrow \infty} \frac{v_N^\pi(s)}{N}$$

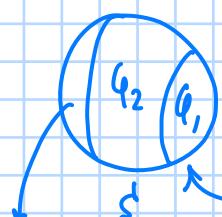
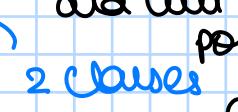
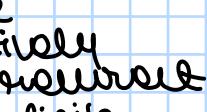
$$g_-^\pi(s) = \lim_{N \rightarrow \infty} \frac{v_N^\pi(s)}{N}$$

Note: $\sup_{\pi \in \Pi^{\text{HR}}} \bar{v}^\pi(s) = \sup_{\pi \in \Pi^{\text{MR}}} \bar{v}^\pi(s)$, instead of Π^{HR} we use Π^{MR}

Finite state MC:

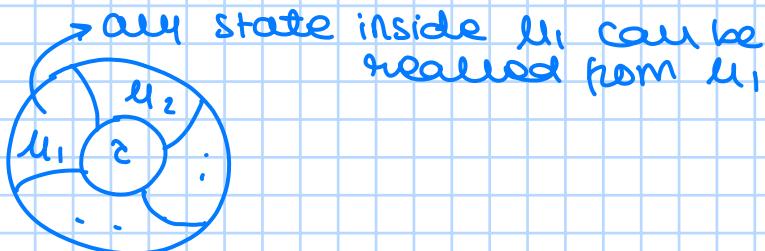
$d^\infty \xrightarrow{\text{for every } d^\infty \text{ we will have one markov chain}}$
 $d^\infty \rightarrow \text{SMR}$

$$P_d(s'|s) = \sum_a P(s'|s, a) q_{d(s, a)} \quad (s| < \infty) \quad \text{finite states}$$

as finite states, there is atleast one recurrent state and will be
   

finitely many classes

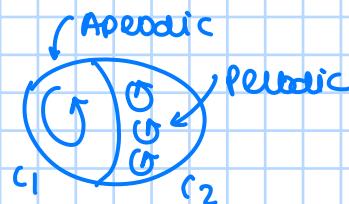
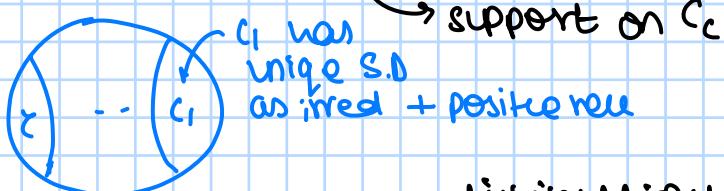
transient class (states we stop visiting after some time)



so, $s \in C$, $\exists \sigma_{s,C} \rightarrow$ prob of getting absorbed in C

= prob of hitting class C before others

Note: Now we will assume stationary distribution exist for every class as S.D exist let $(\mu^{(c)})$ for $s \in C$, $\text{support}(\mu^{(c)}) = \{s : \mu^{(c)}(s) > 0\} = C_c$, $\forall c$



now if S.D exist, $\mu_d P_d = \mu_d$, so $\mu_d^c P_d = \mu_d^c \forall c \in \text{class}$

also, $P_{s,s'}^n = P(X_n = s' | X_0 = s)$ \leftarrow prob of going from $s \rightarrow s'$ in n steps

if M is aperiodic \leftarrow periodicity = 1
 then

$$P_{s,s'}^n \rightarrow \mu^{(c)}(s')$$

if M is periodic then above not guaranteed but

$$\frac{\sum_{k \leq n} p_{s,s'}^k}{n} \xrightarrow{a.s.} \mu^{(c)}(s') \text{ (converges a.s.)}$$

Note: $d^\infty \rightarrow \text{SMR}$ if $s' \in T$ true

$$\frac{\sum_{k \leq n} p_{s,s'}^k}{n} \rightarrow 0 \text{ as } n \rightarrow \infty \text{ in fact } p_{s,s'}^k \rightarrow 0$$

Note: $d^\infty \rightarrow \text{SMR}$, if $s' \in C_c$, $\mu^{(c)} + \text{unac}$, &c

$$\text{then } \frac{\sum_{k \leq n} p_{s,s'}^k}{n} \xrightarrow{} \begin{cases} 0 & ; s \in C_c \text{ with } c \neq c \\ \sigma_{s,c} \mu^{(c)}(s') & ; s \in T \\ \mu^{(c)}(s') & ; s \in C_c \end{cases}$$

now, we are interested in $g_+^\pi(s)$, $g_-^\pi(s)$

$$\begin{aligned} \frac{\mathbb{E}_s \left[\sum_{n=1}^N g_t(x_n) \right]}{N} &= \frac{\sum_{n=1}^N \mathbb{E}_s [g_t(x_n)]}{N} \\ &= \frac{\sum_{n=1}^N \sum_{s'} g_t(s') P(X_n = s' | s)}{N} \\ &= \sum_{s'} g_t(s') \frac{\sum_{n=1}^N p_{s,s'}^n}{N} \xrightarrow{} \begin{cases} \sum_{s' \in C_c} g_t(s') \mu_c(s') & ; s \in C_c \\ \sum_{s \in T} \sigma_{s,c} \mu_c & ; s \in T \end{cases} \end{aligned}$$

$\downarrow \mu_c g_t$

Note: $\mu_c f \stackrel{\Delta}{=} \sum_s \mu_c(s) f(s) = \mathbb{E}[f(X^\infty)]$ where $\sum_s \mu_c(s) = 1$

\checkmark stationary, limit always exist

$$\text{now, } g_+^\pi(s) = \lim_{N \rightarrow \infty} \frac{\mathbb{E}_s^\pi \left[\sum_{n=1}^N g_t(x_n) \right]}{N}$$

$$= \begin{cases} \mu_c \mu & ; s \in C_c \\ \sum_c \sigma_{s,c} \mu_c & ; s \in T \end{cases}$$

μ_c is a distribution

Note: So if only C_c and no T , then $g_+^\pi(s)$ does not depend on s

Theorem: s is finite states, +SMR policy $d^\infty = \pi$, $g_+^\pi = g_-^\pi$ and limit exist

proof: proof follows from above results directly

① unicolor assumption, $\forall d^\infty, \exists$ single positive reward class



② multicolor assumption, here we have multiple classes

Note: under ①, $g^\pi(s)$ will be constant $\forall s$, as $\sigma_{s,c} = 1 \quad \forall s \in S$ as only one c

now, for ①, $\forall s$, $g^\pi(s)$ is same, so we want to try writing D equation

10th Oct:

No class on Tuesday, Tuesday 5:30 class
Next to next Tuesday 5:30 Quiz-2

Quiz-2: 1 or 2 questions, LP solving question

we know $g^\pi(s) = \lim_{N \rightarrow \infty} \frac{V_N^\pi(s)}{N}$ where $V_N^\pi(s) = \sum_{n \leq N} \mathbb{E}_s^\pi [r(x_n, y_n)]$

also, $\underline{g}^\pi(s) = \lim_{N \rightarrow \infty} \frac{V_N^\pi(s)}{N}$

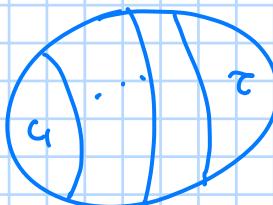
$\overline{g}^\pi(s) = \lim_{N \rightarrow \infty} \frac{V_N^\pi(s)}{N}$

if $d^\infty \rightarrow$ SMR tree $g^\pi(s)$ always exist (finite state space)

π^* optimal if $\underline{g}^{\pi^*}(s) \geq \overline{g}^\pi(s) \forall \pi \rightarrow$ criteria when g^π does not exist

Note: more limits don't exist, we can use above criteria

Finite State models: $\pi = d^\infty$



$\forall \pi = d^\infty$, we have $\{T, c_1, c_2, \dots\}$ clauses
if $s \in C_c^{(\pi)}$, $g^\pi(s) = \mu_c r$ depends on π
depends on π

if $s \in C_c^{(\pi)}$, $g^\pi(s) = \mu_c r = \sum_{s,a} \mu_c(s,a) r(s,a)$

we also saw:

$$d^\infty: P_d \rightarrow M_{d,c}^{(x)} \quad z = (x,y)$$

SMR

$$z_n = (x_n, y_n)$$

$\leftarrow z = (x, z)$ both state and action

prob dist of z_n is $M_{d,c}^{(z)}(s,a) = \underbrace{\mu_{d,c}^{(x)}(s)}_{\text{Prob of state } s} \underbrace{d(s,a)}_{\text{Prob of selecting } a \text{ given } s}$

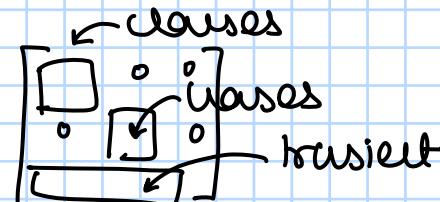
if $s \in T$: $g^\pi(s) = \sum_c \sigma_{s,c} \mu_c r$

in extreme state some π s.t. more space is impossible then only
on μ s.t. $\pi = d^\infty$

$$g^\pi(s) = g \neq s \text{ where } g = Mr$$

Note: as $g^\pi(s)$ does not depend on s , we don't have normal DP equations

for $d^\infty \rightarrow X_n$ under d^∞ : $P_d \sim$



P_d in block matrix form becomes $[]_{n \times m} []_{n \times n}$
 $r \times m \quad \infty \times n - (n+r) \times (n+m)$

If P_d has 2 classes, then $n_1 = |C_1|$, $n_2 = |C_2|$, $q \rightarrow$ transient states
so total $n_1 + n_2 + q$

$$P_d = \begin{bmatrix} P_{d,1} & \overset{n_1 \times n_1}{\underset{0}{\underset{\text{---}}{|}}} & \overset{n_1 \times q}{\underset{0}{\underset{\text{---}}{|}}} \\ \overset{0}{\underset{n_2 \times n_1}{\underset{\text{---}}{|}}} & P_{d,2} & \overset{0}{\underset{n_2 \times q}{\underset{\text{---}}{|}}} \\ \overset{* * * * *}{\underset{q \times (n_1 + n_2 + q)}{|}} & & \end{bmatrix} \quad \text{two need not be zero}$$

$$(n_1 + n_2 + q) \times (n_1 + n_2 + q)$$

now, we can do jordan-canonical form

$$P \chi_i^0 = \lambda_i^0 \chi_i^0$$

$$\begin{bmatrix} \dots \chi_1 \dots \\ \dots \chi_2 \dots \\ \vdots \\ \dots \chi_n \dots \end{bmatrix} \quad W^T P W = J = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

W

if some eigenvalues are repeating, then

$$\begin{bmatrix} \begin{pmatrix} \lambda_1 & 1 \\ & \lambda_1 \end{pmatrix} \\ \begin{pmatrix} \lambda_2 & \\ & \lambda_2 \end{pmatrix} \\ \vdots \\ \begin{pmatrix} \lambda_n & 1 \\ & \lambda_n \end{pmatrix} \end{bmatrix}_{n \times n}$$

The jordan canonical form of P_d :

m eigenvectors of $\lambda = 1$, if $\{c_1, c_2, \dots, c_m, T\}$
+ class, \exists unique stationary distribution

$$\left\{ \begin{array}{l} M_1 = [\underbrace{0 \dots 0}_{n_1 \neq 0} \dots 0] \\ M_2 = [0 \dots 0 \underbrace{0 \dots 0}_{n_2 \neq 0}] \\ \vdots \\ M_m \end{array} \right.$$

Stationary
dist of
m classes

corresponding to each M_i we get: $M_i P_d = M_i^0 \Rightarrow M_i^0$ is eigenvector
of P_d

now, $W = \begin{bmatrix} * \\ M_1 \\ M_2 \\ \vdots \\ M_m \end{bmatrix}$

} any basis (orthogonal by gram-schmidt)
from other eigenvectors
only one vector per class

as stochastic matrix, other eigenvalues < 1

↳ jordan canonical form

$$P_d = W^T \begin{bmatrix} Q & 0 \\ 0 & I \end{bmatrix} W \quad \begin{matrix} m \times m \text{ identity matrix} \\ \downarrow \end{matrix}$$

Note: Now we have $P_d = w^T \begin{bmatrix} Q & 0 \\ 0 & I_{m \times m} \end{bmatrix} w$

now, $\pi_\lambda^\pi = (I - \lambda P_d)^{-1} \pi$ (our discounted case)

but $(I - P_d)^{-1}$ does not exist as eigenvalues 0 (m many)
we use pseudo inverse of $(I - P_d)$

in long run we get gain g if $m=1$, then

$$g^\pi(s) = \mu \pi + s$$

Now we want to know how fast error i.e. $\delta(X_t, Y_t)$
i.e. $\sum_{t \leq N} \mathbb{E} [\pi(X_t, Y_t) - g(X_t)]$ some form of bias and $g(X_t)$ decreases
LLN for markov chains

$$\text{as, } \frac{\sum \mathbb{E} [\pi(X_n, Y_n)]}{N} \rightarrow g(s) \rightarrow \mathbb{E} [\pi(x^\infty, y^\infty)]$$

Pseudo inverse:

$$P = w^T \begin{bmatrix} Q & 0 \\ 0 & I_{m \times m} \end{bmatrix} w \quad \text{now } P^2 = w^T \begin{bmatrix} Q & 0 \\ 0 & w \end{bmatrix}^2 w$$

$$\Rightarrow P^2 = w^T \begin{bmatrix} Q^2 & 0 \\ 0 & I \end{bmatrix} w$$

$$P^k = w^T \begin{bmatrix} Q^k & 0 \\ 0 & I \end{bmatrix} w$$

$$\text{now, } I - P = w^T w - w^T \begin{bmatrix} \theta & 0 \\ 0 & I \end{bmatrix} w$$

$$= w^T \begin{bmatrix} I - \theta & 0 \\ 0 & 0 \end{bmatrix} w$$

Pseudo inverse this is invertible as eigenvalues of θ < 1

$$(I - P)^{\#} = w^T \begin{bmatrix} (I - \theta)^{-1} & 0 \\ 0 & 0 \end{bmatrix} w$$

$$\text{now, } P^* = \lim_{N \rightarrow \infty} P^N \quad \text{if aperiodic}$$

cesaro limit

$$\text{and } P^* = \lim_{N \rightarrow \infty} P^N \quad \text{if periodic}$$

$$= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n \leq N} P^n$$

bounded

$$\text{now, } \frac{\sum \theta^k}{N} = \frac{(I - \theta)^{-1} \theta^N}{N} \rightarrow 0$$

$$\text{so, } P^* = w^T \begin{bmatrix} 0 & 0 \\ 0 & I_{m \times m} \end{bmatrix} w \quad \text{in both cases}$$

Eg: If 2 uses

$n_1 + n_2 + q$ free: $n_1 \rightarrow M_1$ dist

$$P^* = \left[\begin{array}{cccc} \cdots & M_1 & \cdots & \cdots \\ \cdots & M_1 & \cdots & \cdots \\ \cdots & M_1 & \cdots & \cdots \\ \cdots & M_2 & \cdots & \cdots \\ \vdots & & \ddots & \cdots \\ \cdots & M_2 & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \end{array} \right] \left\{ \begin{array}{l} n_1 \text{ times} \\ n_2 \text{ times} \\ q \text{ times} \end{array} \right\} = \left[\begin{array}{ccc} (*)_{n_1 \times n_1} & 0 & 0 \\ 0 & (*)_{n_2 \times n_2} & 0 \\ 0 & 0 & \sigma_{S,1} M_1 + \sigma_{S,2} M_2 \end{array} \right]_{(n_1+n_2+q) \times (n_1+n_2+q)}$$

From above example we get $P^* = \left[\begin{array}{cc} \cdots & M_1 \\ \cdots & M_2 \\ \vdots & \vdots \\ \cdots & M_2 \end{array} \right]$

$$\text{Now, } (I - P + P^*) = W^{-1} \left[\begin{array}{cc} I - \Theta & 0 \\ 0 & 0 \end{array} \right] W \left[\begin{array}{c} \sigma_1 M_1 + \sigma_2 M_2 \\ 0 \\ 0 \end{array} \right] + W^{-1} \left[\begin{array}{cc} 0 & 0 \\ 0 & I \end{array} \right] W$$

$$(I - P + P^*) = W^{-1} \left[\begin{array}{cc} I - \Theta & 0 \\ 0 & I \end{array} \right] W$$

$$(I - P + P^*)^{-1} = W^{-1} \left[\begin{array}{cc} (I - \Theta)^{-1} & 0 \\ 0 & I \end{array} \right] W$$

now we see $PP^* = P^*P = P^*P^* = P$ (from P^* def and we get same)

$$\begin{aligned} PP^* &= \lim_{N \rightarrow \infty} \frac{\sum_{n \leq N} P P^n}{N} \\ \Rightarrow P P^* &= \lim_{N \rightarrow \infty} \frac{\sum_{n \leq N+1} P^n - P}{N} \\ &= P^* \end{aligned}$$

$$\text{Now, } (I - P + P^*)^{-1} - P^* = (I - P)^{\#}$$

$$\begin{aligned} \text{as } (I - P + P^*)P^* &= \cancel{P} - P P^* + P^* P \\ &= P^* \\ \Rightarrow P^* &= (I - P + P^*)^{-1} P^* \end{aligned}$$

$$\text{and so, } (I - P + P^*)^{-1} - P^* = (I - P + P^*)^{-1} (I - P^*) = (I - P)^{\#}$$

Note: we get $(I - P)^{\#} = (I - P + P^*)^{-1} (I - P^*)$ vector of gain

$$\begin{aligned} \text{Now, } (I - P)^{\#} g &= (I - (P - P^*))^{-1} (I - P^*) g \quad \text{where } g = P^* r \\ &= \sum_{k=0}^{\infty} (P - P^*) K (I - P^*) g \quad (\because (I - A)^{-1} = \sum_{k=0}^{\infty} A^k) \end{aligned}$$

$$\begin{aligned} (P - P^*)^2 &= (P - P^*)(P - P^*) \\ &= P^2 - P^* \quad \text{similarly } (P - P^*)K = P K - P^* \end{aligned}$$

$$\text{then, } (\mathbf{I} - \mathbf{P})^{\#} \mathbf{g} = \sum_{k=0}^{\infty} (\mathbf{P}^k - \mathbf{P}^*) (\mathbf{I} - \mathbf{P}^*) \mathbf{g}$$

$$= \sum_{k=0}^{\infty} (\mathbf{P}^k - \mathbf{P}^*) \mathbf{g}$$

$$\text{as for } k=0: (\mathbf{P} - \mathbf{P}^*)^0 (\mathbf{I} - \mathbf{P}^*) \mathbf{g} \\ = (\mathbf{I} - \mathbf{P}^*) \mathbf{g}$$

and $\mathbf{P}^* = \mathbf{P}^k \mathbf{P}^*$ so,

$$(\mathbf{I} - \mathbf{P})^{\#} \mathbf{g} = \sum_{k=0}^{\infty} (\underbrace{\mathbf{P}^k \mathbf{g}}_{\mathbb{E}[r(X_n, Y_n)]} - \underbrace{\mathbf{P}^k g}_{\mathbb{E}[g(X_n)]}) \text{ as } g = \mathbf{P}^* \mathbf{g}$$

$$\text{so, } (\mathbf{I} - \mathbf{P})^{\#} \mathbf{g} = \sum_{k=0}^{\infty} \mathbb{E} [\underbrace{r(X_n, Y_n) - g(X_n)}_{\text{Bias from before}}]$$

Defn: (Bias) Bias will be $(\mathbf{I} - \mathbf{P})^{\#} \mathbf{g}$

Note: for more than one state, we will use CEM and not EM

16th Oct:

$$(I - P + P^*) = W^{-1} \begin{bmatrix} I - Q & 0 \\ 0 & I \end{bmatrix} W$$

$$(I - P)^\# \triangleq W^{-1} \begin{bmatrix} (I - Q)^{-1} & 0 \\ 0 & 0 \end{bmatrix} W$$

$$(I - P + P^*)^{-1} = W^{-1} \begin{bmatrix} (I - Q)^{-1} & 0 \\ 0 & I \end{bmatrix} W$$

$$P^* = W^{-1} \begin{bmatrix} 0 & 0 \\ 0 & I \end{bmatrix} W$$

$$(I - P + P^*)^{-1} - P^* = (I - P)^\#$$

Note: we have $(I - P + P^*)^{-1} P^* = P^*$

$$\text{So, } \underbrace{(I - P + P^*)^{-1} (I - P^*)}_{HP} = (I - P)^\#$$

\leftarrow Pseudo inverse of $(I - P)$

then we get, $\sum (P - P^*)^n (I - P^*) = HP$ (from previous result)

$$\Rightarrow HP = \sum_{n=0}^{\infty} (P^n - P^*) (I - P^*) \quad (\because P^* P = P P^* = P^* P^* = P^*)$$

$$= \sum_{n=0}^{\infty} (P^n - P^*)$$

from λ case: $(I - \lambda P)^{-1} g_i = \vartheta_i$, here we have $(I - P)^\# g_i$
we want to find what comes when

$(I - P)^\# g_i$ is taken

$$\text{now, } \lim_{N \rightarrow \infty} \sum_{n=0}^N (P^n - P^*) g_i = H_P g_i$$

$$\text{now, } \sum_{n=0}^N P^n g_i = \sum_{n=0}^N \mathbb{E}[g_i(x_n, y_n)] \\ = \vartheta_i^\pi$$

$$\text{So, } H_P g_i = \lim_{N \rightarrow \infty} (V_N^\pi - N g^*)$$

$$\text{let } h(s) = \lim_{N \rightarrow \infty} (V_N^\pi(s) - N g^*(s)) = V_N^\pi(s) - N g^*(s) \\ + \sum_{n>N} (P^n - P^*) g_i$$

$$\text{So, } h(s) \approx V_N^\pi(s) - N g_i + O(\downarrow)$$

If we start in
stationarity

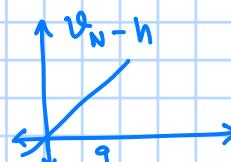
So, pseudo inverse gives us what happens
when we start - start in stationary

$$h(s) - h(j) = V_N^\pi(s) - V_N^\pi(j)$$

$s, j \in \text{same } C_c$

\leftarrow Bias of starting states

really same but fastness in
this argument



periodic:

$$H_{\text{per}} \triangleq \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^N \left(\sum_{k=0}^N p^n - p^N \right) r_k$$

cero sum

or avg sum

$$\frac{1}{N} \sum_{n \in N} v_n(s) - Ng(s) \quad \text{so, } h(s) = \frac{1}{N} \sum v_n(s) - Ng(s) + o(\epsilon)$$

$$\text{so, } h(s) - h(j) = \frac{1}{N} \sum (v_n(s) - v_n(j))$$

Theorem: The gain $g(s) = \lim_{N \rightarrow \infty} \frac{1}{N} v_N(s)$ and $h = H_{\text{per}}$ and $g = (g(s), \forall s)$

(I) true:
(i) $(I - P)g = 0$

(ii) $g + (I - P)h = r$ and

(II) If s_1, s_2, s_3 satisfy (i) and (ii) then they are gain and bias vectors

We assume chain is aperiodic

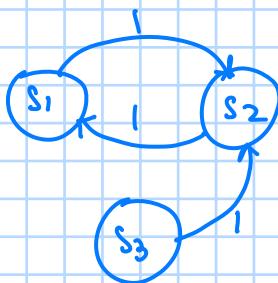
$$\begin{bmatrix} p & 1-p & 0 \\ 1-p & p & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad \{s_1, s_2, s_3\} = S \quad \begin{array}{l} \text{top one is recurrent} \\ \text{bottom is transient} \end{array}$$

If $p=0$: then not aperiodic

$$P>0: \text{ then aperiodic and } P^* = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} \quad \begin{array}{l} \text{Stationary} \\ \text{distribution} \end{array}$$

$$P^N \approx P^* + o(\epsilon)$$

Now if $p=0$:



$$\text{we will have } P^* = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$\text{but } P^N = \begin{cases} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, & N \text{ even} \\ \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, & N \text{ odd} \end{cases} \quad \begin{array}{l} \text{Periodic} \\ \text{Markov} \\ \text{chain} \end{array}$$

$$\lim_{N \rightarrow \infty} P^N = P^*$$

Average and not individual case

$$\text{so in both cases, } P^N \approx \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} + o(\epsilon)$$

Theorem: The gain $g(s) = \lim_{N \rightarrow \infty} \frac{1}{N} v_N(s)$ and $h = H_{\text{per}}$ and $g = (g(s), \forall s)$

(I) true:
(i) $(I - P)g = 0$

(ii) $g + (I - P)h = r$ and

② If g_1, h_1 satisfy (i) and (ii) then they are gain and bias vectors

proof:

$$\textcircled{1} \quad g = P^* h$$

$$(I - P)g = (I - P)P^* h \\ = (P^* - P P^*)h \\ = 0$$

$$\text{Now, } H_P r = (I - P + P^*)^{-1} (I - P^*) = \sum_{n=0}^{\infty} (P^n - P^*)$$

$$(I - P) H_P = H_P (I - P) \quad (\because P P^* = P^* P = P^* P^* = P^*)$$

$$= \sum (P^n - P^*) (I - P)$$

$$= \sum P^n - P^{n+1} - P^{n+1} + P^*$$

$$= \sum P^n - P^{n+1}$$

$$= \lim_{N \rightarrow \infty} \sum_{n=0}^N P^n - \sum_{n=0}^N P^{n+1}$$

$$= \lim_{N \rightarrow \infty} I - P^{N+1} \quad (\because \text{aperiodic})$$

$$= I - P^* \quad \leftarrow \begin{matrix} \text{for periodic case we take claim to get} \\ \text{claim } \lim_{N \rightarrow \infty} P^N = P^* \end{matrix}$$

$$\text{So, } \lim_{N \rightarrow \infty} \left(\frac{1}{N} \sum_{n \in N} (I - P^{N+1}) \right) = I - \lim_{N \rightarrow \infty} P^{N+1} \quad \leftarrow \text{periodic case} \\ = I - P^*$$

$$(I - P) H_P = I - P^* \quad \text{--- ①}$$

$$\text{So, } (I - P) H_P r = I r - P^* r$$

$$\Rightarrow g + (I - P) H_P r = r$$

② g_1, h_1 tuat satisfy (i), (ii)
given

$$(I - P)g_1 = 0$$

and $g_1 + (I - P)h_1 = r$

$$P^* g_1 + P^* \cancel{(I - P)} h_1 = P^* r = g$$

$$\text{So, } (I - P)g_1 + g = g$$

$$\Rightarrow (I - P)g_1 + P^* g_1 = g \quad (\because \text{from ①})$$

$$\Rightarrow (I - P + P^*)g_1 = g$$

$$\Rightarrow (I - P + P^*)g_1 = P^* r$$

$$\Rightarrow \cancel{(I - P + P^*)} g_1 = \cancel{(I - P + P^*)} P^* r$$

$$(\because (I - P + P^*) P^* = P^*)$$

$$\Rightarrow g_1 = P^* r = g$$

$$\Rightarrow g_1 = g$$

$$g_1 + (I - P) h_1 = r$$

$$\Rightarrow P^* r + (I - P) h_1 = r$$

$$\Rightarrow (I - P) h_1 = (I - P^*) r$$

$$\Rightarrow (I - P) h_1 = (I - P) H_P r$$

as $(I - P) H_P = I - P^*$
from ①

$$\Rightarrow (I - P)(h_1 - h) = 0$$

$$\Rightarrow h_1 - h \in N(I - P)$$

$$\Rightarrow h_1 = h + v \text{ for some } v \in N(I - P)$$

Note: Numerically calculating h , we will get different h 's but

$$h - \tilde{h} \in N(I - P)$$

Now, $(I - P)g = 0$ and $g + (I - P)h = r$ and for λ (as we had

$$(I - \lambda P)v = r$$

$$v = \max \{r + \lambda P u\} \rightarrow DP \text{ equation}$$

↑ gain + bias

$$\text{we have } g + h = r + Ph$$

$$\text{Let } i \rightarrow f_i \rightarrow \max_x f_i(x)$$

$$\text{pick set } A_1 \text{ and then do } \max_{x \in A_1} f_2(x)$$

17th Oct:

Theorem: If g and h are gain and bias then

$$(i) \quad (I - P)g = 0$$

$$(ii) \quad g + (I - P)^T h = r$$

and if (g, h) satisfy (i) and (ii) then $g = \text{gain}$
 $h = \text{bias}$

$\lambda \in [0, 1]$ if $\lambda \rightarrow 1$ then we will have average cost version

$$IE[\vartheta] = \frac{1}{T} \sum_{t=1}^T \lambda^{t-1} g_t(s, a)$$

if $\lambda \rightarrow 1$ and $T \rightarrow \infty$ we want to find $IE[\vartheta]$

$$\text{now, } \vartheta_\lambda = \underbrace{(I - \lambda P)^T g}_\text{for } \lambda = 1 \text{ not invertible, so we will take pseudoinverse}$$

for $\lambda \in (0, 1)$ $\lambda = \frac{1}{1+\rho} \quad \rho > 0$

$$= (1+\rho)^{-1}$$

$$\text{then, } (I - \lambda P)^{-1} = (I - (1+\rho)^{-1} P)^{-1} = (1+\rho) \underbrace{(I - P + \rho I)^{-1}}_\text{Result of } (1+\rho)$$

$$R_P = (I - P + \rho I)^{-1}$$

$$I - P = W^{-1} \begin{bmatrix} I - Q & 0 \\ 0 & 0 \end{bmatrix} W$$

$$\begin{aligned} I - P + \rho I &= W^{-1} \begin{bmatrix} I - Q & 0 \\ 0 & D \end{bmatrix} W + W^{-1} \rho I W \\ &= W^{-1} \begin{bmatrix} I - Q + \rho I & 0 \\ 0 & \rho I \end{bmatrix} W \end{aligned}$$

let $I - Q = B \Rightarrow R_P^{-1} = (I - P + \rho I) = W^{-1} \begin{bmatrix} B + \rho I & 0 \\ 0 & \rho I \end{bmatrix} W$

now

$$R_P = W^{-1} \begin{bmatrix} (B + \rho I)^{-1} & 0 \\ 0 & \rho^{-1} I \end{bmatrix} W$$

$$= W^{-1} \begin{bmatrix} (B + \rho I)^{-1} & 0 \\ 0 & 0 \end{bmatrix} W + \frac{1}{\rho} W^{-1} \begin{bmatrix} 0 & 0 \\ 0 & I \end{bmatrix} W$$

$$= W^{-1} \begin{bmatrix} (B + \rho I)^{-1} & 0 \\ 0 & 0 \end{bmatrix} W + \frac{1}{\rho} P^*$$

$$(B + pI)^{-1} = (I + pB^{-1})^{-1} B^{-1}$$

$$= \sum_{n=0}^{\infty} (-1)^n p^n B^{-n}$$

so, $R_p = W^{-1} \begin{bmatrix} \sum_{n=0}^{\infty} (-1)^n p^n B^{-n} & 0 \\ 0 & 0 \end{bmatrix} W + \frac{p^k}{p}$

$$\Rightarrow \vartheta_\lambda = (I - \lambda R_p^{-1}) g = (1 + p) R_p g = (1 + p) p^{-1} p^k g$$

$$+ (1 + p) W^{-1} \begin{bmatrix} \sum_{n=0}^{\infty} (-1)^n p^n (I - Q)^{-n} & 0 \\ 0 & 0 \end{bmatrix} W g$$

$$H_p = W^{-1} \begin{bmatrix} (I - Q)^{-1} & 0 \\ 0 & 0 \end{bmatrix} W$$

$$\Rightarrow \vartheta_\lambda = (1 + p) p^{-1} p^k g + (1 + p) \sum_{n=0}^{\infty} (-1)^n p^n H_p^{n+1} g$$

$$\lambda = (1 + p)^{-1}$$

as $p \rightarrow 0$

so, $\vartheta_\lambda = \left(1 + \frac{1}{p}\right) p^k g + \underbrace{\sum_{n=1}^{\infty} (-1)^n p^n H_p^{n+1} g}_{= H_p g} + \underbrace{\sum_{n=1}^{\infty} (-1)^n p^n H_p^{n+1} g}_{= H_p g}$

$$\Rightarrow \vartheta_\lambda = \left(1 + \frac{1}{p}\right) p^k g + h + p \left[\underbrace{\sum_{n=1}^{\infty} (-1)^n p^n H_p^{n+1} g}_{f(\lambda) \text{ term}} + \underbrace{\sum_{n=0}^{\infty} (-1)^n p^n H_p^{n+1} g}_{\text{this is bounded}} \right]$$

$f(\lambda) \rightarrow 0$ as $\lambda \rightarrow 1$

$$\Rightarrow \vartheta_\lambda = (1 - \lambda)^{-1} g + h + f(\lambda)$$

as $\lambda \rightarrow 1$, $h \rightarrow 0$, $f(\lambda) \rightarrow 0$

$$\text{so, } g = \lim_{\lambda \rightarrow 1} (1 - \lambda) \vartheta_\lambda$$

this is just an heuristic to understand when $\lambda \rightarrow 1$ we get

If $\Sigma \sim \text{geom } \sigma \cdot v$

$$g = \vartheta_\lambda$$

$$\begin{aligned} \mathbb{E} [\vartheta_\lambda(X_\Sigma)] &= \sum_n P(\Sigma = n) \mathbb{E} [\vartheta_\lambda(X_n)] \\ &= \sum_{n=0}^{\infty} (1 - \lambda) \lambda^{n-1} \mathbb{E} [\vartheta_\lambda(X_n)] \\ &\quad \uparrow \text{failing with prob } \lambda \\ &= (1 - \lambda) \sum_{n=0}^{\infty} \lambda^n \mathbb{E} [\vartheta_\lambda(X_n)] \end{aligned}$$

$$\Rightarrow \mathbb{E}[r(X_\lambda)] = (1-\lambda)\mathcal{V}_\lambda$$

as $\lambda \rightarrow 1$
 $\zeta \rightarrow \infty$ almost surely
 and so,

$$g = \lim_{\lambda \rightarrow 1} (1-\lambda)\mathcal{V}_\lambda$$

one more way to look what happens
 when $\lambda \rightarrow 1$

optimality equation:

$$\max \{ r_d + (\lambda P_d - I) \mathcal{V}_\lambda^* \} = 0$$

$$\text{or } \mathcal{V}_\lambda^* = \max \{ r_d + \lambda P_d \mathcal{V}_\lambda^* \}$$

$$\text{now } \mathcal{V}_\lambda = (1-\lambda)^T g + h + f(\lambda)$$

$$\Rightarrow \max \left\{ r_d + (\lambda P_d - I) [(1-\lambda)^T g + h + f(\lambda)] \right\} = 0$$

if $\lambda \rightarrow 1$ then

Note: If we assume unique, g becomes g.e, $g \in \mathbb{R}$
 $r_d + (\lambda P_d - I) [(1-\lambda)^T g + h + f(\lambda)]$

$$= r_d + (\lambda P_d - I)(1-\lambda)^T g.e + (\lambda P_d - I)h$$

$\xrightarrow{g \in \mathbb{R}}$ all one vector

$$= r_d + \lambda P_d \cdot e g (1-\lambda)^T - (1-\lambda)^T g e + (\lambda P_d - I)h$$

$$= r_d + \lambda g.e (1-\lambda)^T - (1-\lambda)^T g e + (\lambda P_d - I)h$$

$$= r_d - g.e (1-\lambda)^T (1-\lambda) + (\lambda P_d - I)h$$

$$= r_d - g.e + (\lambda P_d - I)h$$

$$0 = r_d - g.e + (\lambda P_d - I)h$$

so, optimality equation becomes:

$$\sup_d \{ r_d - g.e + (\lambda P_d - I)h \} = 0$$

Theorem: If $\max_d \{ r_d - g.e + (\lambda P_d - I)h \} \leq 0$ for some (g, h) then
 $g \cdot e \geq g^*$

Now from above theorem, we can see if we take min of
 all such g 's we will get g^*

re min g

$$g + h(s) - \sum_j p(j|s,a) n(j) \geq r(s,a) \quad \forall s,a \quad \left. \right\} \text{LP for average cost mdp}$$

then dual becomes

$$\max_{\pi(s,a)} \sum_{s,a} r(s,a) \pi(s,a)$$

s.t.

$$\sum_a \pi(j,a) - \sum_{s,a} p(j|s,a) \pi(s,a) = 0$$

$$\sum_{s,a} \pi(s,a) = 1$$

and $\pi(s,a) \geq 0 \forall s,a$

$\left. \right\}$ Dual of mdp

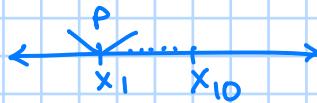
POMDP:

POMDP are Partially Observable MDPs, a particular case of POMDP is IOMDP or intermittently observable MDP

IOMDP:

t X_t is observable with probability p
is lost with $1-p$

$$S = \{s_1, \dots, s_n\}$$



b_t is a vector

$$b_t = P(s_t | s_1, \dots, s_{t-1}, a_1, \dots, a_{t-1}, a_t)$$

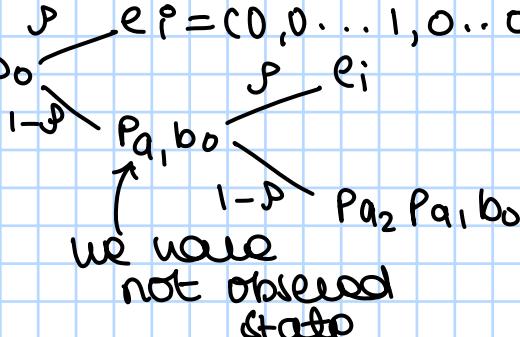
only ones
which are observed

$$= P(s_t | \underbrace{s_{t-k}, a_{t-k}, a_{t-k+1}, \dots, a_t}_{\text{last observed state}})$$

All actions taken

after last observed state

then $b_0 \xrightarrow{p} e^0 = (0, 0 \dots 1, 0 \dots 0)$ as we have observed new state



transition probabilities

$$Q(b_{t+1} | b_t, a_t) = \begin{cases} p [P_{a_t} b_t]_0 & ; b_{t+1} = e^0, i \in S \\ (1-p) & ; b_{t+1} = P_{a_t} b_t \\ 0 & ; \text{otherwise} \end{cases}$$

we have observed s_{t+1}

$P_{a_t} b_t$

we have not observed b_{t+1}

Note: b^0 is a vector s.t. $b = (b(s), \forall s)$

now, for IOMDP, two states are b_t 's

$$\text{so, } R(b, q) = \sum_s b(s) M(s, a) \leftarrow \text{expected value}$$

Average reward IOMDP: function of s as b is a vector

we want $\frac{1}{T} \sum_{t=1}^T \mathbb{E}[R(b_t, a_t)]$ as $t \rightarrow \infty$

then we have an avg reward MDP with new state space which is not countable

$$\text{dual: } \max_{\pi(b, a)} \pi(b, a) R(b, a)$$

$$\text{s.t. } \sum_a \pi(b', a) - \sum_{b, a} Q(b' | b, a) \pi(b, a) = 0 \quad \forall b'$$

$$\sum_{b, a} \pi(b, a) = 1 \text{ and } \pi(b, a) \geq 0 \quad \forall b, a$$

as $b \in [0, 1]$ then uncountable many states and we can't solve MDP

but: ① If underlying MDP has finite states and action set

② underlying MDP is recurrent

under ①, ② we have strong duality and we can solve MDP

thus makes b countable infinite space

21st Oct:

we know $(I - P_d)g = 0$, $n - g + (P_d - I)h = 0$

we have, $\pi = d^\infty$

$n = n_d$
 $g = g_d$

If g and h satisfy the equations above and $P_d^* h = 0$
 then $h = H_p h \rightarrow$ Bias

optimality:

$$g_+^\pi = \limsup_{N \rightarrow \infty} \frac{1}{N} \mathcal{V}_N^\pi \quad \pi \in \Pi^{HR}$$

↑ any history dependent random policy

$$g_-^\pi = \liminf_{N \rightarrow \infty} \frac{1}{N} \mathcal{V}_N^\pi$$

Defn: (optimality) π^* is optimal if $g_-^{\pi^*} \geq g_+^{\pi^*} + \pi$

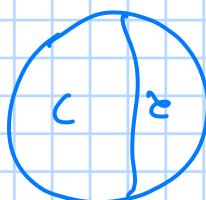
Note: $g_-^\pi \leq g_+^\pi + \pi$ and equality if optimal

$$\text{Defn: } g_+^* = \sup_{\pi \in \Pi^{HR}} g_+^\pi$$

$$g_-^* = \sup_{\pi \in \Pi^{HR}} g_-^\pi$$

uniqueness:

If $d^\infty \in MR$ policy, the chain has almost sure the recurrent class



$$g_d^\infty = \underbrace{g \cdot e}_{\text{all one vector}} \quad g \in \mathbb{R}$$

also we don't care about periodic/aperiodic as we take limit

$$\text{let } g_d = g_d \cdot e \quad e = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad d^\infty \in MR$$

$$\Rightarrow \underbrace{(I - P_d)g_d}_{} = 0$$

$\underbrace{(I - P_d)(e)}_{e - e = 0} = 0$
 this is meaningless

$$n_d - g_d e + (P_d - I)h_d = 0 \quad (\because \text{second equation})$$

$$\mathcal{V}_\lambda^{d^\infty} = (I - \lambda)^{-1} g_d e + h + f(\lambda) \quad (\because \text{from previous result})$$

↑ converge

$\exists \lambda < 1$ s.t. $\mathcal{V}_\lambda^* = \mathcal{V}_\lambda^{d^\infty}, \forall \lambda > \bar{\lambda}$ differently for different d^∞
 so, optimal policy = $\delta^\infty, \forall \lambda > \bar{\lambda}$

We also know, $\forall d^\infty, \lim_{\lambda \rightarrow 1} v_\lambda^{d^\infty}(1-\lambda) = g^{d^\infty}$

the optimal policy g^∞ remains same for every
 $\xleftarrow{\text{if } g^\infty \text{ is optimal}}$ $\xrightarrow{\text{any cost}} \uparrow$
 $\xleftarrow{\lambda}$ $\xrightarrow{\text{optimal}}$ $\xrightarrow{\text{if } \lambda \text{ suff big}}$
 $\xleftarrow{\text{scalar vector}}$

We did max to find

$$v_\lambda^* = \max_d \{ r_d + \lambda p_d v_\lambda^* \}$$

$$\text{let } B(g, h) \triangleq \max_d \{ r_d - g_e + (p_d - I)h \}$$

then our optimality equations says if $B(g, h) = 0$ then g is optimal gain

- Theorem : (i) If $\exists g$ scalar and vector h s.t $B(g, h) \leq 0$ then $g \geq g_+$
(ii) If $\exists g$ scalar and vector h s.t $B(g, h) \geq 0$, then $g \leq g_-$
(iii) If $\exists g, h$ s.t $B(g, h) = 0$ then $g_e = g^* = g_- = g_+$
 $\&$, $\exists d^\infty$ s.t $g = g_d = g^*$

Proof : (i) Given a vector and scalar (h, g) s.t
 $B(g, h) \leq 0$

Let $\pi = (d_1, d_2, d_3, \dots) \in \Pi^{MR}$

$\forall d, r_d - g_e + (p_d - I)h \leq 0 \quad \forall d \in \Pi^{MR}$

$\Rightarrow g_e \geq r_d + (p_d - I)h \quad \forall t = 1, 2, \dots$ (for π)

We want to show $g_e \geq g_+^\pi \quad \forall \pi \in \Pi^{MR}$ (it implies for all HR as we can map HR and MR)
as $\forall \pi \in \Pi^{HR}, \exists \tilde{\pi} \in \Pi^{MR}$ s.t $g_+^\pi = g_+^{\tilde{\pi}}$
 $g_-^\pi = g_-^{\tilde{\pi}}$

Now, $g_e = p_d, g_e \geq p_{d_1} (r_{d_2} + (p_{d_2} - I)h)$ from ① as p_d is positive

similarly $p_{d_1} p_{d_2} g_e$

$\vdots \geq p_{d_1} p_{d_2} (r_{d_3} + (p_{d_3} - I)h)$

$\forall t \quad p_{d_1} p_{d_2} \dots p_{d_t} g_e$

$\geq p_{d_1} \dots p_{d_t} (r_{d_{t+1}} + (p_{d_{t+1}} - I)h)$

Summing this $\forall t \leq N$

↳ telescopic

$$\Rightarrow N g_e \geq \sum_t p_{\pi, t} r_{d_{t+1}} + \sum_t (p_{\pi, t+1} - p_{\pi, t}) h$$

where $p_{\pi, t} = p_{d_1} p_{d_2} \dots p_{d_t}$

$$\Rightarrow N g_e \geq \sum_t p_{\pi, t} r_{d_{t+1}} + (p_{\pi, N+1} - I)h$$

$$\Rightarrow g \leq \frac{1}{N} \sum_t^N \pi_t^* r_{t+1} + \frac{1}{N} (\pi^*, N+1) h \xrightarrow[N \rightarrow \infty]{\text{Bounded}} 0$$

$$\Rightarrow g \leq \underline{g}, \underline{g}^{\pi} \geq \underline{g}_{\pi} \quad \forall \pi \in \Pi^{MR}$$

$$\Rightarrow g \leq \underline{g}^{\pi} \quad \forall \pi \in \Pi^{MR} \text{ and } \underline{g} \leq \underline{g}^{\pi} \quad \forall \pi \in \Pi^{HR}$$

(ii) given a vector and scalar (h, g) , s.t $B(g, h) \geq 0$

$$\Rightarrow B(g, h) = \max_a (\dots) = g_d^* - g e + (P_d^* - I) h \geq 0$$

$$\Rightarrow g \leq g_d^* + (P_d^* - I) h$$

$$\text{then } P_d^* g \leq P_d^* (g_d^* + (P_d^* - I) h)$$

$$\Rightarrow g \leq g_d^* e \quad (\text{from previous part})$$

$$\Rightarrow g \leq g^* e \quad (\because \text{By defn of } g^*)$$

given (i), (ii), part 3 is trivial and follows

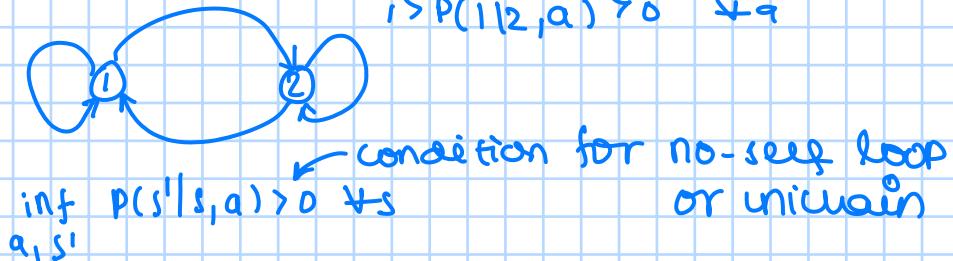
$$g^* \leq g \leq g_d^* \leq g^* \leq g^*$$

$$\Rightarrow g = g^* = g^* = g_d^*$$

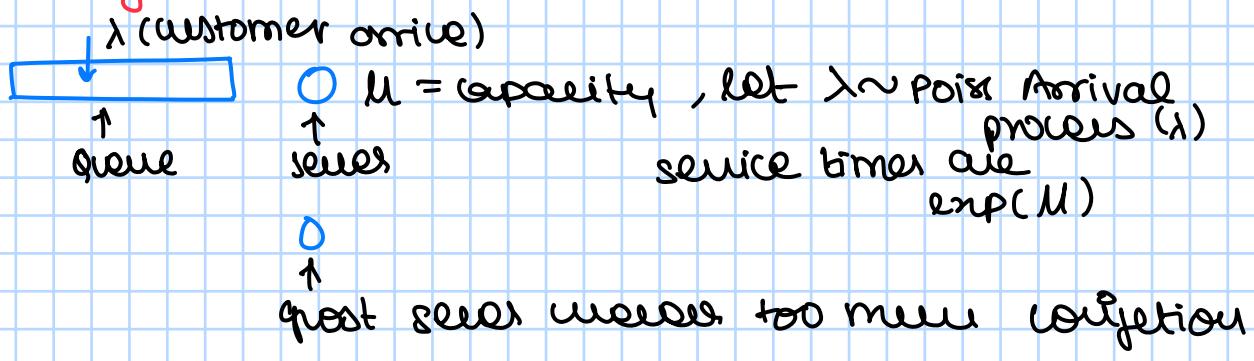
eg: $A_1 = \{a_{11}, a_{12}\}$ $A_2 = \{a_{21}, a_{22}\}$ $g_1(s, a)$

$$1 > P(1|1, a) \geq 0 \quad \forall a$$

$$1 > P(1|2, a) \geq 0 \quad \forall a$$



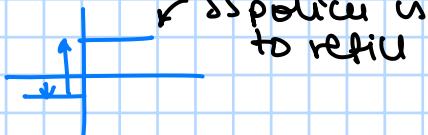
Queuing model:



C_s = additional server cost, policy: $\exists \alpha \text{ s.t. } \forall s < 0 \quad \alpha^s = 0$
otherwise $\alpha^s = 1$

Incentive control:

$(w, c_s) \quad c_0 = k + \tilde{c}_0(q) \quad \text{if } k > 0 \quad \text{ss policy:}$



24 Oct:

we have $B(g, h) = 0$, $g \rightarrow \text{scalar}$, $h \rightarrow \text{vector}$

$$B(g, h) = \max_d \{ r_d + (P_d - I)h - g e \}$$

we have known that if
 $B(g, h) \leq 0 \Rightarrow g \geq g^*$
 $B(g, h) > 0 \Rightarrow g \leq g^+$
 $B(g, h) = 0 \Rightarrow g = g^*$

Exercise:

Theorem: If $\|\cdot\| \leq \infty$, then \exists scalar g , vector h s.t $B(g, h) = 0$
assuming inchain

we have already seen, $\exists \bar{\lambda} < 1$ s.t $\forall \lambda > \bar{\lambda}$, the unique optimal policy $d^* = \delta + \lambda > \bar{\lambda}$ is same

Theorem: If $\|\cdot\| \leq \infty$, then \exists scalar g , vector h s.t $B(g, h) = 0$
assuming inchain

Proof: consider $\lambda \uparrow 1$ let $v_{\lambda n} \rightarrow$ value vector $\forall \lambda n$

$$v_{\lambda n} = \max_d \{ r_d + \lambda_n P_d v_{\lambda n} \}$$

$$\Rightarrow 0 = \max_d \{ r_d + (\lambda_n P_d - I) v_{\lambda n} \}$$

fixing a $v_{\lambda n}^{d^\infty}$
 $v_{\lambda n} = (I - \lambda_n)^{-1} g + h + f(\lambda)$ where $v_{\lambda n} = v_{\lambda n}^{d^\infty}$
comes from δ^∞ , so does h

$$0 \geq r_d + (\lambda_n P_d - I)(I - \lambda_n)^{-1} g$$

$$g \leq g^{d^\infty}, h = H_p r_g$$

$$0 \geq r_d + (\lambda_n P_d - I)((I - \lambda_n)^{-1} g + h + f(\lambda_n))$$

$$\Rightarrow 0 \geq r_d + (P_d - I)h + (\lambda_n - 1)P_d h - g \leq + f(\lambda_n) \quad (.)$$

as $\lambda n \rightarrow 1$, h , g does not change as δ^∞ policy

as $\lambda n \uparrow 1$

$$0 \geq r_d + (P_d - I)h - g \leq \forall d$$

for $d = \delta$ the \geq converts into equality as it's the maximiser

$$\text{so, } 0 = r_\delta + (P_\delta - I)h^{d^\infty} - g^{d^\infty} e \Leftrightarrow B(g^{d^\infty}, h^{d^\infty}) = 0$$

Two cases Age of Information:

we have 2 sources, when they update age=0 if updated
so, we are assuming some initial condition

age \rightarrow age+1 if not updated

State space: $X = (x_1, x_2)$ age of source 1, age of source 2

now, if $x_1 \neq x_2$ (Stealing) tree

$$\{ = \{ (x_1, x_2) \mid x_1 \neq x_2 \}$$

$$= \{ \}_1 \cup \{ \}_2$$

$$\{ \}_1 = \{ (x_1, x_2) \mid x_1 < x_2 \}$$

$$\{ \}_2 = \{ (x_1, x_2) \mid x_2 < x_1 \}$$

↑ as if we do not
start at equal
we can never get to
it

we need a smaller state space if possible

(i) it is still markov

(ii) no loss of options

we can make $\{ \}'_1 = \{ (x_1, x_2) \mid x_1 < x_2 \text{ s.t. } x_2 \leq L \}$

$$\{ \}'_2 = \{ (x_1, x_2) \mid x_2 < x_1 \text{ s.t. } \begin{array}{l} x_2 \leq L-1 \\ x_1 \leq L \end{array} \}$$

if we see same policy s.t.
in $\{ \}'_1 \rightarrow q_1 = 1, \{ \}'_2 \rightarrow q_2 = 2$

then this policy is not unique
as 2 different immediate clauses

Algorithm for Unique models :

$$V^0 \leftarrow \arg \max_d (\pi_d + \lambda P_d V^0)$$

$$V^1 \leftarrow d_1$$

⋮

$$V^n \leftarrow \arg \max_d (\pi_d + \lambda P_d V^{n-1})$$

$$V^n \leftarrow d_n$$

Above is for λ or discounted cost MDP, as $\lambda < 1$ we got contraction

$$\left\{ \begin{array}{l} L_d V = \pi_d + \lambda P_d V \quad \forall d \text{ linear for every } d \\ d V = \max_d L_d V \quad \text{non-linear} \end{array} \right.$$

Both are contractions

$$|d V - d V'| \leq \lambda |V - V'|$$

$$V^{n+1} = f(V^n) \rightarrow V^*$$

↑ unique fixed point

Aug Cost :

$$0 = \max_d \{ \pi_d + (\lambda P_d - I) V \}$$

$$0 = \max_d \{ \pi_d + (\lambda P_d - I) h - g^T \}$$

↑ do not want this

we want V here

Defn: (Span norm) Not a norm, but

$$\text{sp}(\mathbf{v}) = \max_i v_i - \min_i v_i$$

If we take $\text{sp}(\mathbf{v})$, $\text{sp}(\mathbf{v} + g\mathbf{e}_i)$

we have $\text{sp}(\mathbf{v}) = \text{sp}(\mathbf{v} + g\mathbf{e}_i)$

thus solves our issue with $g\mathbf{e}_i$ vector

now, $|f(\mathbf{v}) - f(\mathbf{v}')| \leq \lambda |\mathbf{v} - \mathbf{v}'|$ is what we want

$$\text{sp} T\mathbf{v} \leq \alpha \text{sp}(\mathbf{v}) \quad \text{for } \alpha < 1$$

↑
linear
operator

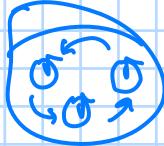
then we want to see if we will get solution

$$\text{sp}(T^n \mathbf{v}) \leq \alpha^n \text{sp}(\mathbf{v}) \rightarrow 0$$

now, if we know d^* , then

$$\begin{aligned} Td^*(n) &= g_{d^*} + P_{d^*} h - g\mathbf{e}_i \\ Td^*(n-n') &= \cancel{g_{d^*}} + P_{d^*}(h-h') - \cancel{g_{d^*}} \\ &= P_{d^*}(h-h') \end{aligned}$$

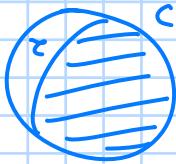
23rd Oct:



unichain if periodic chain $P^n \rightarrow p^*$

i.e. $p^{dn+j} \rightarrow p^*$
as $n \rightarrow \infty$, j/d is the period

Multichain models:

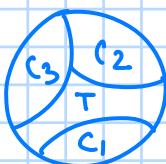


The multichain has only one recurrent class and so

$$\frac{g_d}{\uparrow} = g_d^e \leftarrow \text{constant/same gain vector}$$

$\frac{\sum_{n \leq N} (P^{d^\infty})_s^n}{N} \rightarrow g$ as $N \rightarrow \infty$, given $|f| < \infty$ or finite state space

Note: If $\limsup_{T \rightarrow \infty} \frac{1}{T} |f| = \infty$, then we can also have cases where all



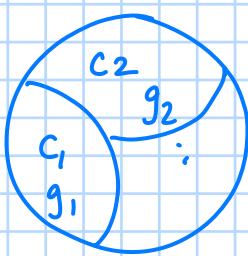
the markov chain will depend on d

Ex: \exists unique recurrent class $\forall d^\infty, d \in \pi^{\text{SMD}}$
iff \exists unique recurrent class $\forall d^\infty, d \in \pi^{\text{SMR}} \rightarrow$ done down

Note: for unichain from above we can only take $\pi \in \pi^{\text{SMD}}$

Defn: (multichain) If $\exists d \in \text{MD}$, with more than one recurrent class, then it is multichain

Properties of g_d :



$$\begin{pmatrix} g_1 \\ g_2 \\ \vdots \\ g_2 \\ g_3 \\ \vdots \\ g_3 \end{pmatrix}$$

$= g$, same for every $s \in$ same recurrent class

Policy evaluation eqn:

given d^∞ ,

(i) $(P_d - I)g = 0$

(ii) $\pi_d + P_d h - h - g = 0$

{ this works for both unichain and multichain

we have (gain, bias) vectors satisfy above, and if

some (g', h') satisfy above equation, then
 $g' = g$ and $h' = h + \text{Null}(I - P_d)$

Note: If $P_d h = 0$ then $h' = h$, bias vector is h

optimality equations:

$$\text{we need } \max_{d \in M} (P_d - I) g = 0$$

$$\text{and } \max_{d \in M} (q_d + P_d h - h - g) = 0$$

if (h, g) satisfy above, but we will have an issue
as many diff d's satisfy diff equation

$$\text{so, } \max_d (P_d - I) g = 0$$

$$\max_{d \in E} (q_d + P_d h - h - g) = 0 \text{ where }$$

$$E = \{d : d(s) = a, \text{ where } (P_d - I) g = 0\}$$

so, g should be s.t. $(P_d - I) g \leq 0 \quad \forall d \in M$

$$\sum_{j \in S} p(j|s, a) g(j) - g(s) \leq 0 \quad \forall s, a \quad \text{--- (1)}$$

and \exists at least one a' for fixed s s.t (1) is satisfied with equality

$$\sum_{j \in S} p(j|s, a') g(j) - g(s) = 0$$

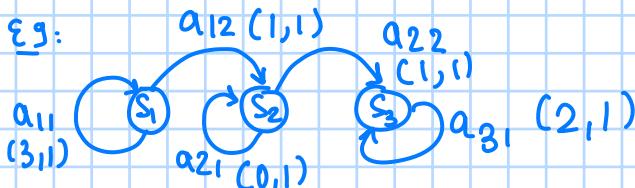
$$\text{Hence, let } B_s = \{a' \in A_s \mid \sum_{j \in S} p(j|s, a') - g(s) = 0\} \quad \text{--- (2)}$$

$$\max_{a' \in B_s} (q(s, a') + \sum_j p(j|s, a') h(j) - h(s) - g(s)) = 0, \forall s \quad \text{--- (3)}$$

$$\text{so, } E = \{d \mid d(s) \in B_s \quad \forall s\}$$

Theorem: If (g^*, h^*) vectors satisfy multichain OE and are bounded, then $\exists M < \infty$ s.t. $\exists (g^*, h^* + Mg^*)$ satisfying modified multichain OE

e.g:



$$p(S_1 | S_1, a_{11}) = 1 \quad p(S_2 | S_1, a_{12}) = 1 \quad p(S_2 | S_2, a_{21}) = 1 \quad p(S_3 | S_2, a_{22}) = 1$$

$$q(S_1, a_{11}) = 3$$

$$q(S_1, a_{12}) = 1$$

$$q(S_2, a_{21}) = 0$$

$$q(S_2, a_{22}) = 1$$

$$q(S_3, a_{31}) = 1$$

$$q(S_3, a_{31}) = 2$$

from Unicolumn Optimality equation, g is a scalar and

$$\max_d (r_d + P_d h - h - g e) = 0$$

$$\text{i.e. } g + h_1 = \max \{ 3 + h_1, 1 + h_2 \}$$

$$h_1^* = h(s_1)$$

$$g + h_2 = \max \{ 0 + h_2, 1 + h_3 \}$$

$$g + h_3 = \max \{ 2 + h_3 \} = 2 + h_3$$

if we assume unicolumn

$$g = 2 \text{ from last}$$

$$2 + h_2 = \max \{ h_2, 1 + h_3 \}$$

$$2 + h_1 = \max \{ 3 + h_1, 1 + h_2 \}$$

$$3 + h_1$$

\Rightarrow this is a contradiction

optimal policy: $g^*(s_1) = 3$ if on s_1 we stay on s_1 ,

$g^*(s_2) = g^*(s_3) = 2$ if on s_2 we go to s_3 and stay same for s_3

Theorem: If \exists a pair of vectors (g, h) that simultaneously satisfy

$$\max_d (P_d - I) g \leq 0$$

$$\text{and } \max_d (r_d + P_d h - h - g) \leq 0$$

$$\text{then (i) } g \geq g^* = \sup_{\pi \in \Pi^{HR}} \limsup_{N \rightarrow \infty} \frac{\mathbb{V}_{N+1}^\pi}{N+1}$$

if $\exists (g, h)$ for some d' s.t. $(P_{d'} - I)g > 0$ and $r_{d'} + P_{d'} h - h - g > 0$

$$\text{then (ii) } g \leq g^{(d')} \leq g^* = \sup_{\pi \in \Pi^{HR}} \liminf_{N \rightarrow \infty} \frac{\mathbb{V}_{N+1}^\pi}{N+1}$$

if $\exists (g, h)$ satisfying $\max_d (P_d - I)g = 0$ and $\max_d (r_d + P_d h - h - g) = 0$

$$\text{then (iii) } g = g^*$$

with below theorem, we have multicolumn and modified OE

Theorem: If (g^*, h^*) vectors satisfy multicolumn OE and are bounded, then $\exists M < \infty$ s.t. $\exists (g^*, h^* + Mg^*)$ satisfying modified multicolumn OE

Proof: Let (g^*, h^*) satisfy multicolumn optimality equations

say for some (s, a)

$$g = r(s, a) + \sum_j p(j|s, a) h^*(j) - h^*(s) - g^*(s) > 0$$

$\Rightarrow a \notin B_s$ as g^*, h^* satisfy multicolumn OE

Consider $h = h^* + Mg^*$ for some $M > 0$

$$\begin{aligned}
q(s, a) + \sum_j p(j|s, a) h(s) - h(s) - g^*(s) \\
&= q(s, a) + \sum_j p(j|s, a) h^*(s) - h^*(s) - g^*(s) \\
&\quad + M(\sum p(j|s, a) g^*(j) - g^*(s)) \\
&= c_1 + M(\sum p(j|s, a) g^*(j) - g^*(s)) \\
&\quad \text{as } a \notin B_S \underbrace{\sum p(j|s, a) g^*(j) - g^*(s)}_{c_2} < 0
\end{aligned}$$

then for $M = -\frac{c_1}{c_2}$
as $c_1 > 0, c_2 < 0 \Rightarrow M > 0$

or $(g^*, h^* + Mg^*)$ satisfy modified OE

Note: From above theorem, we have both modified OE and multi-chain OE

Sol: \exists unique recurrent class $\forall d^\infty, d \in \pi^{\text{SMD}}$
iff \exists unique recurrent class $\forall d^\infty, d \in \pi^{\text{SMR}}$

Ans: If \exists unique recurrent class, $\forall d^\infty \in \pi^{\text{SMR}}$

then as $\pi^{\text{SMD}} \subseteq \pi^{\text{SMR}}, \forall d^\infty \in \pi^{\text{SMD}}$
 \exists unique recurrent class

now, \exists unique recurrent class $\forall d^\infty \in \pi^{\text{SMD}}$
true

let $d : \mathcal{S} \rightarrow P(A)$ be st $d^\infty \in \pi^{\text{SMR}}$

now if \nexists unique recurrent class for $d^\infty \in \pi^{\text{SMR}}$
true let

T, C_1, C_2, \dots, C_k be classes $\kappa \gamma_1 \Rightarrow \kappa \gamma_2$
true

$g_d \neq g \in \mathbb{R}^\infty, \exists s_1, s_2 \text{ s.t. } g_d(s_1) \neq g_d(s_2)$

$$g_d(s_1) = \sum q_d(s_1, a) g_d(s_1, a) \xrightarrow{\text{same}} = \sum q_d(s_2, a) g_d(s_2, a)$$

$= g_d(s_2)$ this is a contradiction

31st Oct :

multichain MDPs:

$$\max_d (P_d - I) g = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{multichain OE}$$

$$\max_{d \in E} r_d + (P_d - I) h - g = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{if } d \in MD$$

\Updownarrow Finite MDPs

$$\max_d (P_d - I) g = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{modified OE} \quad \leftarrow \begin{array}{l} \text{if } (g^*, h^*) \text{ satisfy multi OE} \\ \text{then } (g^*, h^* + Mg^*) \text{ satisfy modi OE} \end{array}$$

$$\max_d r_d + (P_d - I) h - g = 0$$

where $B_S = \{a \in A_S \mid \sum p(j|s, a) g(j) - g(s) = 0\} \neq \emptyset$

$$E = \{d \mid d(s) \in B_S, \forall s\}$$

Theorem: If \exists a pair of vectors (g, h) that simultaneously satisfy

$$\max_d (P_d - I) g \leq 0$$

$$\text{and } \max_d (r_d + P_d h - h - g) \leq 0$$

then (i) $g \geq g_+^* = \sup_{\pi \in \Pi^{HR}} \limsup_{N \rightarrow \infty} \frac{\mathbb{V}_N^\pi}{N+1}$

if $\exists (g, h)$ for some d' s.t. $(P_{d'} - I) g > 0$ and $r_{d'} + P_{d'} h - h - g > 0$

then (ii) $g \leq g_-^{(d')} \leq g_-^* = \sup_{\pi \in \Pi^{HR}} \liminf_{N \rightarrow \infty} \frac{\mathbb{V}_N^\pi}{N+1}$

if $\exists (g, h)$ satisfying $\max_d (P_d - I) g = 0$ and $\max_d (r_d + P_d h - h - g) = 0$

then (iii) $g = g^*$

proof: (i) we have $\exists (g, h)$ s.t. $\max_d (P_d - I) g \leq 0$, $\max_d (r_d + (P_d - I) h) \leq 0$
 $\Rightarrow (P_d - I) g \leq 0$ and $r_d + (P_d - I) h - g \leq 0 \quad \forall d \in MD$

now $g_+^* = \sup_{\pi \in \Pi^{HR}} \lim_{N \rightarrow \infty} \frac{\mathbb{V}_N^\pi}{N+1}$

take any $d \in MD$, then

$$(P_d - I) g \leq 0, \quad r_d + (P_d - I) h - g \leq 0$$

let $\pi \in \Pi^{HR}$, let $\pi = (d_1, d_2, \dots) \rightarrow r_{d_1} + (P_{d_1} - I) h \leq g$

$$g \geq P_d g \Rightarrow g \geq P_d g \quad \forall i \quad (\because \text{①})$$

then $g \geq P_{d_2} g \geq P_{d_1} (r_{d_2} + (P_{d_2} - I) h)$

$$\vdots$$

$$g \geq P_{d_n} \dots P_{d_1} g \geq P_{d_n} \dots P_{d_1} (r_{d_{n+1}} + (P_{d_{n+1}} - I) h)$$

$$\Rightarrow Ng \geq \sum_t p\pi_{1t} r_{d_{t+1}} + \sum_t (p\pi_{1t+1} - p\pi_{1t}) h \quad p\pi_{1t} = P_{d_1} P_{d_2} \dots P_{d_t}$$

$$\geq \sum_t p\pi_{1t} r_{d_{t+1}} + (p\pi_{1,n+1} - I) h$$

$$\Rightarrow g \geq \frac{1}{N} \sum_t p\pi_t \gamma d_{t+1} + \underbrace{\frac{1}{N} (Pd - I) h}_{\substack{\text{Bounded} \\ \rightarrow 0 \text{ as } N \rightarrow \infty}}$$

$$\Rightarrow g \geq g_+^\pi \forall \pi$$

$$\Rightarrow g \geq \sup_\pi g_+^\pi = g_+^*$$

(ii) $\pi_{d^*} = g + (Pd^* - I) h \geq 0$ and $(Pd^* - I) g \geq 0$
for some (g, h, d^*)

$$\Rightarrow g \leq \pi_{d^*} + (Pd^* - I) h \text{ and } g \leq Pd^* g$$

$$\Rightarrow Pd^* g \leq Pd^*(\pi_{d^*} + (Pd^* - I) h)$$

$$\Rightarrow g \leq Pd^*(\pi_{d^*} + (Pd^* - I) h)$$

then $g \leq \frac{1}{N} \sum_t p\pi_t \gamma \pi_{d^*} = g_-^{d^*}$

$$\Rightarrow g \leq g_-^*$$

(iii) from (i), (ii) this follows trivially

Note: In λ case, v_λ^* was unique but policy not unique, also MD policies are finite

we know that if $\lambda_n \rightarrow 1$, then we can take a subsequence s.t g^∞ is optimal & $\lambda_n \rightarrow$ subsequence,

$\exists \lambda^* < 1$ s.t g^∞ optimal & $\lambda \geq \lambda^*$, g^∞ is optimiser for average cost MDP

Theorem: (Existence) If $|f| < \infty$, $|h| < \infty$, $\forall \pi \in \mathcal{P}$, then \exists a solution to multichain OE

Proof: $g = g^{s^\infty}, h = h^{s^\infty}$ are true candidates

$$0 = \pi_g + (\lambda_n P_g - I) v_{\lambda_n} \geq \pi_d + (\lambda_n P_d - I) v_{\lambda_n} \quad \forall d \in MD$$

(\because optimality of $\lambda < 1$ case)

$$0 \geq \pi_d + (\lambda_n P_d - I) \left(\frac{g}{(1-\lambda_n)} + h + f(\lambda_n) \right) \quad \forall d \in MD$$

$$= \pi_d + \frac{(P_d - I) g}{1-\lambda_n} - \frac{(\lambda_n - 1) P_d g}{1-\lambda_n} + (P_d - I) h + (\lambda_n - 1) P_d h + \tilde{f}(\lambda_n)$$

$$= \frac{(P_d - I) g}{1-\lambda_n} + \underbrace{(\pi_d - P_d g + (P_d - I) h)}_{\substack{\lambda_n \rightarrow 1 \\ \text{true term remains}}} + \underbrace{(\lambda_n - 1) P_d h + \tilde{f}(\lambda_n)}_{\substack{\rightarrow 0 \\ \text{as } \lambda_n \rightarrow 1}}$$

as $\lambda_n \rightarrow 1$
if $(P_d - I) g < 0$
then not $\rightarrow \infty$

Case I: $(P_d - I) g > 0$ is not possible as after some λ_n we cannot satisfy the equation

Case II: $(P_d - I)g = 0$

$$\text{then } 0 \geq g_d - P_d g + (P_d - I)h \neq d$$

$$\Rightarrow 0 \geq g_d - g + (P_d - I)h \text{ as } (P_d - I)g = 0$$

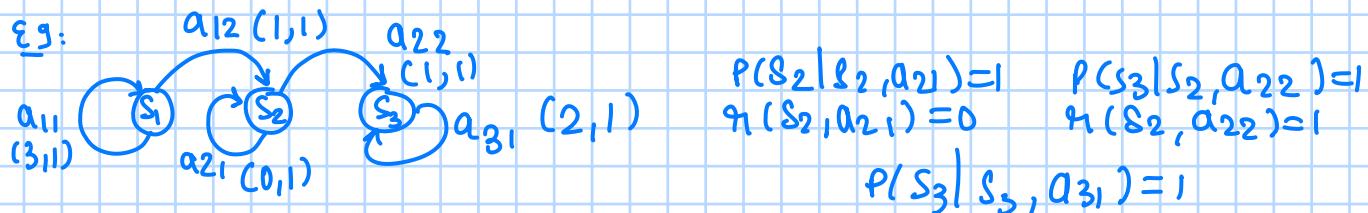
Case III: $(P_d - I)g < 0$, this case is possible, but d does not matter as in multilinear OE, we will only take $d: (P_d - I)g = 0$ from Case II and Case III

$$\max_d (P_d - I)g = 0 \quad (\because d = \infty, \text{ we have } P_{\infty}g = g)$$

$$\max_{d \in E} (g_d + (P_d - I)h - g) = 0 \quad (\because d = \infty \text{ we have } \max = 0)$$

Note: Above we have enstace of multilinear OE, we will have enstace of modified OE by $h^* + Mg^*$ theorem

Eg:



$$P(S_2 | S_2, a_{21}) = 1 \quad P(S_3 | S_2, a_{22}) = 1$$

$$g_1(S_2, a_{21}) = 0 \quad g_1(S_2, a_{22}) = 1$$

$$P(S_3 | S_3, a_{31}) = 1$$

$$P(S_1 | S_1, a_{11}) = 1 \quad P(S_2 | S_1, a_{12}) = 1$$

$$g_1(S_1, a_{11}) = 3 \quad g_1(S_1, a_{12}) = 1$$

$$g_1(S_3, a_{31}) = 2$$

$$g_1, g_2, g_3 \\ h_1, h_2, h_3$$

$$\max\{3 + h_1, 1 + h_2\} = h_1 + g_1$$

$$\max\{0 + h_2, 1 + h_3\} = h_2 + g_2 \quad \left. \begin{array}{l} \text{second multilinear} \\ \text{OE} \end{array} \right\}$$

$$\text{and } 2 + h_3 = h_3 + g_3$$

$$\left. \begin{array}{l} g_1 = \max\{g_1, g_2\} \\ g_2 = \max\{g_2, g_3\} \end{array} \right\} \text{first multilinear OE}$$

$$g_3 = g_3$$

$$\text{as } g_1 > 3 \Rightarrow g_1 > g_3 \text{ and } g_1 > g_2, g_3 = 2$$

$$\text{if: } g_1 = g_2$$

but as $\max\{0 + h_2, 1 + h_3\} = h_2 + g_2$
we pick a_{22} then $g_2 = g_3 \Rightarrow \text{contradiction}$

$$\text{so, } g_1 > g_2 \Rightarrow g_1 > 3, g_2 = 2 = g_3 \text{ then } g_1 = 3$$

as if $3 + h_1 < 1 + h_2$ then we pick a_{12}
but $g_1 > 3$ and not $g_1 = g_2 = 2$
so contradiction

$$\text{so } 3 + h_1 = g_1 + h_1 \Rightarrow g_1 = 3$$

$$\text{so } g^* = (3, 2, 2)^T$$

Lemma: If (g^*, h^*) satisfy multi OE and $d^* \in \operatorname{argmax}_d \{ \mathbb{r}_d + P_d h^* \}$
 then $(d^*)^\infty$ is optimal

Lemma: If (\tilde{g}, \tilde{h}) satisfies modified OE and d^* s.t. $P_{d^*} \tilde{g} = \tilde{g}$ and
 $d^* \in \operatorname{argmax}_d \{ \mathbb{r}_d + P_d \tilde{h} \}$
 then $(d^*)^\infty$ is optimal

policy improvement algorithm:

for unichain: d_0

Step n: find g_n, h_n

$$\mathbb{r}_{d_n} + (P_{d_n} - I) h_n - g_n e = 0$$

improvement: $d_{n+1} \in \operatorname{argmax}_d (\mathbb{r}_d + (P_d h_n))$

if $d_{n+1} = d_n$ we stop
 first preference of d_{n+1} is d_n , obtain as
 many $s | d_n(s) = d_{n+1}(s)$ possible

for multichain: do

$$\text{step } n: \mathbb{r}_{d_n} + (P_d - I) h_n - g_n = 0 \\ (P_{d_n} - I) g_n = 0$$

improve: $d_{n+1} \in \operatorname{argmax}_d (P_d g_n)$

