

IE 708

Markov Decision Processes

Ref: M. Puterman
Markov decision process
grading: class test 1 - 15%
(before midterm)
class test 2 - 20%.
midterm - 30%.
final exam/presentation
pre req: James Norris - chapter 1
- 35%.
read

29th July:

There is a notion of time, at a given time we have to make a decision, take them s.t. we maximise something or some function.

Eg: Running competition, speed fast or less at start (speed to end), control speed at each time, finish race first

Very broadly about sequential decision making under uncertainty

Eg: first step, we want to keep stock, goal is to maximise profit, every morning order something s.t. we meet demand and also less stock for next day doesn't want to order too many, minimise wastage, so how much to buy?

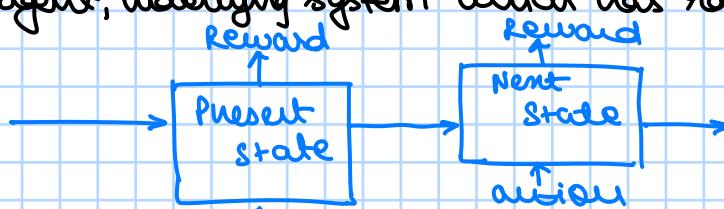
Running objectives:

- 1) Model sequential decision problems using MDPs. Problems
- 2) Existence and structure of optimal strategies
- 3) Computing an optimal strategy

Markov
Decision
Problems

Framework:

There is a single agent, underlying system which has randomness, at a given time:



↑ Action ← this action can be chosen in many ways and
single agent we assume agent knows everything about past

← time →

The decision maker wants to optimise action to maximise some form of reward
reward can be MA or total average

↑ Moving average /
over period of time

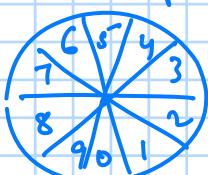
Eg: For vegetable vendor:

State = stock available
action = how much to order
reward = revenue

Two sequential problems:

Who's calling:

↓ Disperse some numbers



First let's have only one team, after every spin we place num at some location,
goal is to max out the num.

In optimal strategy, current state does not matter, only location matters

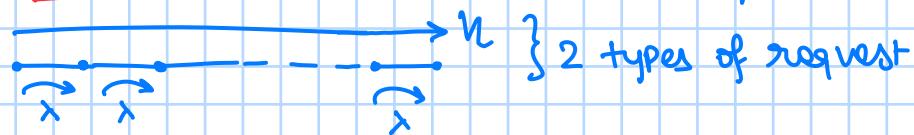
Obs Num spin 1 spin 2 Spin 3 Spin 2 spin 1

0	5	4	3	2	1
1	5	4	3	2	1
2	5	4	3	2	1
3	4	3	3	2	1
4	3	3	2	2	1
5	3	2	2	1	1
6	2	2	1	1	1
7	1	1	1	1	1
8	1	1	1	1	1
9	1	1	1	1	1

Average score using this strategy: 78,734.12

Average score using random assignemt: 49,999.5

K-hop tandem network: (Resource allocation problem)



each link has a unit capacity

reward of 1-hop request: 1

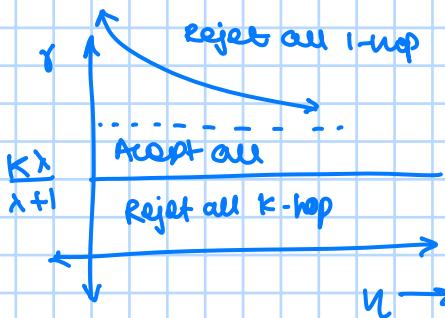
reward of K-hop request: τ

On an average tree are λ arrivals (for one unit of time)

K -hop request arrive at N

parameters: λ, N, τ, K

↑ reward of K-hop
↑ per unit time K-hop req



Model description & Notation:

1) Decision epoch:

$T \rightarrow$ set of epoch, can be discrete time / random time

most of the time

for finite discrete time, i.e. $T = \{1, 2, \dots, N\}$
we call it finite horizon problem

Random time is when the period
is random

for infinite discrete time, i.e. $T = \{1, 2, \dots\}$ we call it infinite horizon problem

2) State and actions:

S : set of states

mostly discrete / finite
countable

A_s : set of allowable actions for present state $s \in \mathcal{S}$

$A = \bigcup_{s \in \mathcal{S}} A_s$: set of all actions

mostly focus on discrete

Note: $P(A_s)$ = set of probability distributions on A_s

$q(\cdot) \in P(A_s)$ a random action, and so

$\therefore q(a)$ for $a \in A_s$ is a random feasible action

e.g.: $A_s = \{1, 2, \dots, a_s\}$

↑
set of possible actions

$$P(A_s) = \left\{ q \mid \sum_{i=1}^{a_s} q_i = 1, q_i \geq 0 \right\}$$

$q(\cdot) \in P(A_s)$

↖ A vector of probability and so

$q(a)$ for $a \in A_s$ is a random thing from A_s set

3) Reward and state evaluation:

if in state $s \in \mathcal{S}$ and we take action $a \in A_s$
then reward we get is $r_t(s, a)$

$r_t(s, t)$: reward at time t , when system is at $s \in \mathcal{S}$ and we take action $a \in A_s$

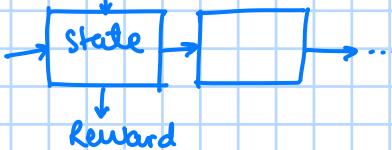
$p_t(j|s, a)$: the prob that system is at $j \in \mathcal{S}$ given prev is $s \in \mathcal{S}$ and we take $a \in A_s$ at given time t . ↙ this is on $t+1$

Note: $r_t(s, a, j)$ is reward if we are on $s \in \mathcal{S}$, take $a \in A_s$ and we go to $j \in \mathcal{S}$
in this case:

$$r_t(s, a) = \sum_{j \in \mathcal{S}} p_t(j|s, a) \times r_t(s, a, j)$$

Defn: $(T, \mathcal{S}, A_s, s \in \mathcal{S}, r_t(s, a), p_t(j|s, a))$ is the markov decision process

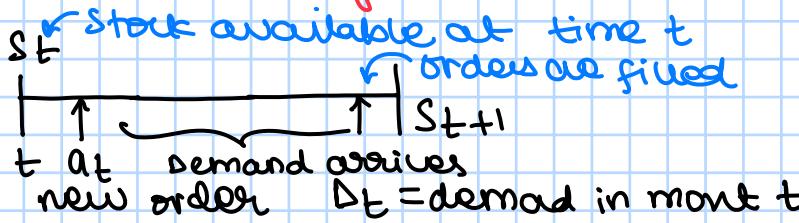
1st Aug:



$(\mathcal{T}, \mathcal{S}, \mathcal{A}, s \in \mathcal{S}, r_t(s, a), p_t(j|s, a))$ is called a MDP

$T = \{1, 2, 3, \dots, N\}$ random horizon (stop when we find candidate)
 ↗ one of the cases (N is a random variable)
 decision epochs are random is one more case of decision epochs

Stochastic inventory control:



$$P(\Delta_t = j) = p_j$$

M = total capacity of warehouse so, $a_t + s_t \leq M$

$\mathcal{T} = \{1, 2, 3, \dots, N\}$ $N < \infty$ for finite horizon
 $N = \infty$ for infinite horizon

$\mathcal{S} = \{0, 1, \dots, M\} \rightarrow$ state / items in warehouse

$\mathcal{A}_S = \{0, 1, \dots, M-S\} \quad s \in \mathcal{S}$
 ↗ maximum we can order

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

$$s_{t+1} = (s_t + a_t - \Delta_t)^+ \leftarrow \text{one way of waiting}$$

$h(u)$ = holding u units of stock per month
 $o(u)$ = ordering u units

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

some fixed cost
cost of u

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

if current state s_t and we order a_t units then

$$\text{holding cost} = h(s_t + a_t)$$

$$\text{ordering cost} = o(a_t)$$

$$\text{revenue cost} = f(\max \{D_t, s_t + a_t\})$$

$$\pi_t(s_t, a_t) = f(\max \{D_t, s_t + a_t\}) - o(a_t) - h(s_t + a_t)$$

$$\gamma_t(s_t, a_t, s_{t+1}) = f(s_t + a_t - s_{t+1}) - o(a_t) - h(s_t + a_t)$$

\downarrow Avg profit if we start with u units

$$F(u) = \sum_{j=0}^u f(j) P(j) + f(u) [\sum_{j>u} P(j)] \rightarrow \text{expected profit for given } u \text{ units}$$

$$\tau_t(s_t, a_t) = F(s_t + a_t) - h(s_t + a_t) - \sigma(a_t)$$

$$\text{or } \tau_t(s, a) = F(s+a) - h(s+a) - \sigma(a)$$

$$F(u) = \mathbb{E} [f(\min_{t'} D_t, u)] \quad D_t \text{ is a random variable}$$

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

Ques: So who's courting game MDP formulation? \rightarrow doubt

Decision rules:

prescribes an action a_t at decision epoch t

$$d_t: S \rightarrow A$$

so $d_t(s) \in A_s$ action when current state is s
also called Markovian Deterministic (MD)

$$(A_s \subseteq A \text{ so } d_t(s) \in A_s \subseteq A)$$

$h_t = (s_1, a_1, s_2, a_2, \dots, s_{t-1}, a_{t-1}, s_t)$ is called history upto time t

$$h_t \in \underbrace{S \times A \times \dots \times S}_{1} \times \underbrace{A \times \dots \times S}_{t-1} \times \underbrace{S}_t$$

$$d_t: H_t \rightarrow A$$

$d_t(h_t) \in A_{S_t}$ History dependent deterministic (HD)

M. Randomised $d_t: S \rightarrow P(A)$

H. Randomised $d_t: H_t \rightarrow P(A)$

D_t^K : set of all decision rules of type K

$$K \in \{MD, HD, MR, HR\}$$

policy:

Specifies the decision rule at each decision epoch

$$\Pi = (d_1, d_2, \dots)$$

$\Pi = (d, d, \dots)$ \rightarrow this is called a stationary

Π^K = set of all policies for $K \in \{SD, SR, MD, MR, HD, HR\}$ policy

$$\Pi^{SD} \subseteq \Pi^{MD} \subseteq \Pi^{MR} \subseteq \Pi^{HR}$$

$$\Pi^{SD} \subseteq \Pi^{MD} \subseteq \Pi^{HD} \subseteq \Pi^{HR}$$

$$\Pi^{SR} \times \Pi^{MD}$$

$$\Pi^{HD} \times \Pi^{MR}$$

$$\text{Ex: } d(s) = \begin{cases} 0 & \text{if } s > 0 \\ M-s & \text{if } s < 0 \end{cases}$$

↓
fixed threshold
 $d : \{ \rightarrow A$ and as
 $\pi = (d_1, d_2, \dots)$
we have

$\pi \in \Pi^{SD}$

if

$$d_t(s) = \begin{cases} 0 & ; s > t \\ M-s & ; s \leq t \end{cases}$$

then $d_t : \{ \rightarrow A$
and $\pi = (d_1, d_2, \dots)$
 $\Rightarrow \pi \in \Pi^{NS}$

5th Aug:

$\Omega = \underbrace{S \times A \times S \times A \times \dots \times S \times A}_{\text{canonical space}} \Omega^N$ Finite horizon ($N < \infty$)

also assume finite states^N and actions^N

$B(\Omega) = \text{set of all subsets of } \Omega$ (this will be a σ -field as finite Ω)
 $\omega \in \Omega$

$$\omega = (s_1, a_1, s_2, a_2, \dots, s_{N-1}, a_{N-1}, s_N)$$

$X_t(\omega) = s_t$ state at time t

$Y_t(\omega) = a_t$ action at time t

$Z_t(\omega) = (s_1, a_1, s_2, a_2, \dots, s_t)$ history upto time t

now if $\Pi \in \Pi^{MR}$ (the biggest policy space)

$$\Pi = (d_1, \dots)$$

$$\text{s.t. } d_t: Z_t \rightarrow P(A)$$

this induces a probability distribution P^Π

let $\text{distn}(X_1) = p_1$, so under P^Π :

$P^\Pi(X_1 = s) = p_1(s) \leftarrow \text{probability that at } t=1, \text{ state is } s$

$P^\Pi(Y_t = a | Z_t = h_t) \leftarrow \begin{array}{l} \text{given } \Pi, \text{ what is prob we take} \\ \text{action } a \text{ at time } t \\ \text{given history at time } t \\ Z_t = h_t \end{array}$

here $d_t: H_t \rightarrow P(A)$

$$\text{so } d_t(h_t) \in P(A)$$

$q(a)$ is a probability

now, $P^\Pi(X_{t+1} = s_{t+1} | Z_t = h_t, Y_t = a) = p_{t+1}(s_{t+1} | s_t, a)$

\leftarrow probability of next state is s_{t+1} given history and action

$$P^\Pi(s_1, a_1, s_2, a_2, \dots, s_{N-1}, a_{N-1}, s_N) = p_1(s_1) \times q_{d_1(s_1)}(a_1) \times p_2(s_2 | s_1, a_1) \times q_{d_2(s_2)}(a_2) \times \dots \times p_{N-1}(s_{N-1} | s_{N-2}, a_{N-2}) \times p_N(s_N | s_{N-1}, a_{N-1})$$

so, $(\Omega, B(\Omega), P^\Pi)$ describes the MDP

Note: (X_t, Y_t) will not be markovian here as action will depend on history

P^Π for $\Pi \in \Pi^{MR}$:

$$P^\Pi(X_{t+1} = s_{t+1} | Z_t = h_t)$$

$$= \sum_{a \in A} P^\Pi(X_{t+1} = s_{t+1}, Y_t = a | Z_t = h_t) \quad (\text{law of total prob})$$

$$= \sum_{a \in A} P^\Pi(Y_t = a | Z_t = h_t) \times P^\Pi(X_{t+1} = s_{t+1} | Y_t = a, Z_t = h_t)$$

$$= \sum_{a \in A} q_{d_t(s_t)}(a) p_t(s_{t+1} | s_t, a)$$

\leftarrow this is standard conditional prob

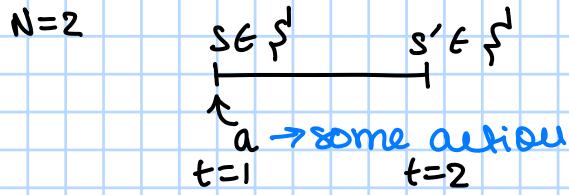
$$= \sum_{a \in A} q_{d_t(s_t)}(a) p_t(s_{t+1} | s_t, a)$$

so, $\{X_t\}$ is a discrete time markov chain (DTMC) by defn

Ex: (X_t) is DTMC, then is $(f(X_t))$ a markov chain? $f: S \rightarrow S'$
Ans: If f is one-one it is DTMC but if not one-one then not a markov chain

Note: $(X_t, r_t(X_t, Y_t))$ is called a markov reward process, not necessarily a MC

one period markov decision problem:



At time 2, we get a reward $r(s')$ for some $r: S' \rightarrow \mathbb{R}$

If we took action a' at $t=1$, then we get reward $r_t(s, a')$

$$\sum_{j \in S} r(j) P_t(j|s, a') = \text{Average reward} = \frac{\mathbb{E}_{\pi}^{\pi}[r(X_2) | X_1=s]}{\mathbb{E}_{\pi}^{\pi}[r(X_2)]}$$

$$\text{total reward} = r_t(s, a') + \sum_{j \in S} r(j) P_t(j|s, a')$$

so given s stable at $t=1$, the action to take is

$$a^* = \underset{\substack{a' \in A \\ \text{MD}}}{\operatorname{argmax}} \left\{ r_t(s, a') + \sum_{j \in S} r(j) P_t(j|s, a') \right\}$$

MR case:

$$d_i, q \in P(A_S)$$

$$\max_a \sum_j q(a) [r_t(s, a) + \sum_j r(j) P(j|s, a)] = \text{maximum reward}$$

$$\text{Note: } \max_{a' \in A_S} \left\{ r_t(s, a') + \sum_j r(j) P_t(j|s, a') \right\} = \max_{q \in P(A_S)} \left\{ \sum_j q(a) [r_t(s, a) + \sum_j r(j) P(j|s, a)] \right\}$$

so in this case randomisation does not help

so if we start at s :

$$\mathbb{E}_S^{\pi} [r_t(X_1, Y_1) + r(X_2)] \quad \begin{array}{l} \text{fix state, this is the avg reward} \\ \text{we get if we are using } \pi \text{ policy} \end{array}$$

Note: For finite case we can find an action to maximise reward, but for infinite case not always true

Sequence of rewards:

$N < \infty$, finite horizon case $R = (R_1, \dots, R_N)$ sequence of rewards
we can compare random variables by:

U, V U stochastically dominate V if
 $P(U > t) \geq P(V > t) \forall t$

now if $U = (U_1, \dots, U_N)$
 $V = (V_1, \dots, V_N)$ then we can compare U, V by:

$(E[f(U_1 \dots U_N)]) \geq E[f(V_1 \dots V_N)]$ for any f s.t
 $0 \leq f \Rightarrow f(U) \leq f(V)$ (this is coordinate wise)

(this is a very strong condition)

Eg: π_1, π_2 we want to compare two policies, $N=2$

$$P^{\pi_1}(R=(0,0)) = P^{\pi_1}(R=(1,0)) = P^{\pi_1}(R=(0,1)) = P^{\pi_1}(R=(1,1)) = \frac{1}{4}$$

$$P^{\pi_2}(R=(0,0)) = P^{\pi_2}(R=(1,1)) = \frac{1}{2}$$

$$f(U,V) = U+V \text{ then } E^{\pi_1}[f] = (2) \times \frac{1}{4} + 1 \times \frac{1}{4} + 1 \times \frac{1}{4} = 1$$

$$E^{\pi_2}[f] = 2 \times \frac{1}{2} = 1$$

by this fusion π_1 and π_2 are same 2

$$f(u,v) = \max\{u,v\}$$

$$E^{\pi_1}[f] = \frac{1}{4} + \frac{1}{4} + \frac{1}{4} = \frac{3}{4}$$

$$E^{\pi_2}[f] = \frac{1}{2} \text{ so } \pi_1 \text{ is better than } \pi_2$$

similarity for $f = \min\{u,v\}$ π_2 is better than π_1

Note: The above example is why we cannot use the method of fusions to compare

Expected total reward $\pi \in \Pi^{HR}$

$r_N(s)$: reward at time N
at time N we get

under policy π $\rightarrow E_S^\pi \left[\sum_{t=1}^{N-1} r_t(x_t, y_t) + r_N(x_N) \right]$ = reward by starting at s and using policy π

Initial state s (start from s) = total expected reward under π , given $x_1=s$

$$= V_N^\pi(s)$$

goal of decision maker is to maximise $V_N^\pi(s)$

or $\sup_{\pi \in \Pi^{HR}} V_N^\pi(s) \rightarrow$ goal of policy maker to pick $\pi \in \Pi^{HD}$ with $\max V_N^\pi(s)$

8th Aug:

Quiz - Aug 29

$X_1 = S$, given $\Pi \in \Pi^{HR}$ policy Π (Random), we get $(\mathcal{S}, \mathcal{B}(\mathcal{S}), P^\Pi)$ where $N < \infty$ (finite horizon problem)
we get $r_t(x_t, y_t)$

$$E_s^\Pi \left[\sum_{t=1}^{N-1} r_t(x_t, y_t) + r_N(x_N) \right] = V_N^\Pi(s) \quad \begin{array}{l} \text{final reward depends} \\ \text{upon final state} \end{array}$$

\uparrow
conditioned
to start at s

Expected total reward starting at $x_1 = s$ under Π
 S, A are discrete

if $\sup_{s,a} |r_t(s, a)| < \infty$, $\sup_s |r_N(s)| < \infty$ then $V_N^\Pi(s)$ is always finite, $V_N^*(s) = \sup_{\Pi \in \Pi^{HR}} V_N^\Pi(s)$ (\Rightarrow \exists supremum)

Defn: (optimal policy) a policy $\Pi^* \in \Pi^{HR}$ is optimal if

$$V_N^*(s) = V_N^{\Pi^*}(s) \geq V_N^\Pi(s) \quad \forall s \in S \quad \begin{array}{l} \text{so optimal} \\ \text{policy does not} \\ \text{always exist} \end{array}$$

(i.e. $V_N^{\Pi^*}(s) = \sup_{\Pi \in \Pi^{HR}} V_N^\Pi(s)$)

Defn: (ϵ optimal) Π is ϵ optimal if

$$V_N^\Pi(s) + \epsilon \geq V_N^*(s) \quad \forall s \in S \quad (\text{this will always exist})$$

Note: From defn of supremum, there will always be an ϵ optimal strategy

$E[V_N^\Pi(X_1)]$ for some strategies also get maximised by Π^*

now, $E_s^\Pi \left[\sum_{t=1}^{N-1} \lambda^{t-1} r_t(x_t, y_t) + \lambda^{N-1} r_N(x_N) \right]$ expected total discount reward

so we want to give more reward to first few time (since use of λ (idle or $d(f)$) ($\lambda < 1$, but not comp))

Note: Expected total discounted reward and normal reward will be maximised similarly

$V_N^\Pi(s)$ for given Π ($\Pi \in \Pi^{HP}$):

Defn: (reward-to-go function) $U_t^\Pi : H_t \rightarrow \mathbb{R}$

$$U_t^\Pi(h_t) = E_{h_t}^\Pi \left[\sum_{n=t}^{N-1} \gamma_n(x_n, y_n) + \gamma_N(x_N) \right]$$

$(h_t = (s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t))$

Expected reward from t to N given history till t h_t

Now at $t=t$ action taken is $d_t(h_t)$ ($d_t(h_t) = a_t \in A_{s_t}$)
and given s_t so

$$U_t^\Pi(h_t) = r_t(s_t, d_t(h_t)) + E_{h_t}^\Pi \left[\sum_{n=t+1}^{N-1} \gamma_n(x_n, y_n) + \gamma_N(x_N) \right]$$

$$\text{now } E[f(X)|X] = f(X)$$

$$E[E[Y|X]] = E[Y] \rightarrow \text{tower property}$$

$\mathcal{F} \supseteq \sigma$ 2 σ -fields

$$\mathbb{E}_{N-1}^{\pi} [\mathbb{E}[X|\mathcal{F}]|h_t] = \mathbb{E}[X|h_t]$$

$$X = \sum_{n=1}^{t+1} r_n(X_n, Y_n) + r_{t+1}(X_{t+1})$$

$$\text{so, } U_t^{\pi}(h_t) = r_t(s_t, d_t(h_t)) + \mathbb{E}_{h_{t+1}}^{\pi} [\mathbb{E}_{h_{t+1}}^{\pi} [X]]$$

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

$$\text{Now, set } t=N, U_N^{\pi}(h_N) = r_N(s_N) + h_N$$

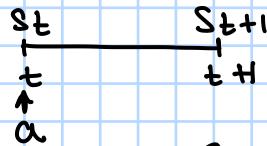
only thing
that is
random

Substitute $t-1$ in place of t to evaluate ①
to compute $V_{t-1}(h_{t-1}) + h_{t-1} \in H_{t-1}$

and then stop when $t=1$, $U_1^{\pi}(s) + s \in S$
and so, $V_N^{\pi}(s) = U_1^{\pi}(s)$

Note: The above algorithm to calculate $V_N^{\pi}(s)$ is called policy evaluation algorithm

optimal reward-to-go:



$$U_t(h_t) = \sup_{a \in A_{s_t}} \left\{ r_t(s_t, a) + \sum_{j \in S} p_t(j | s_t, a) V_{t+1}(h_{t+1}, a, j) \right\} = U_t^*(h_t)$$

not for given π like before

choose action
that maximizes the above

$$U_N(h_N) = r_N(s_N)$$

Note: The above equation is optimality equation/Bellman equation

$$\text{Now, let } U_t^*(h_t) = \sup_{\pi \in \Pi^{\text{HP}}} V_t^{\pi}(h_t)$$

But we still have to show U_t^* will satisfy the Bellman eqn

Theorem: Let U be a solution to optimality equation, then $U_t(h_t) = U_t^*(h_t)$
 $+ h_t \in H_t$, in particular $U_t(s) = V_N^*(s) + s \in S$

Proof:

We will show that $U_n(h_n) \geq U_n^*(h_n) + h_n$ and then we will
show that $\forall \epsilon > 0, \exists \pi \in \Pi^{\text{HP}}$ s.t.

$$U_n^{\pi}(h_n) + (N-n)\epsilon \geq U_n(h_n)$$

$$\text{then, } U_n^{\pi}(h_n) + (N-n)\epsilon \geq U_n(h_n) \geq U_n^*(h_n)$$

but by defn: $U_n^*(h_n) \geq U_n^{\pi}(h_n)$
of sup

$$\Rightarrow U_n^*(h_n) + (N-n)\epsilon \geq U_n(h_n) \geq U_n^*(h_n)$$

and we let $\epsilon \rightarrow 0$

$$U_t(h_t) = \sup_{a \in A} \left\{ r_t(s_t, a) + \sum_{j \in S} p_t(j|s_t, a) \cdot U_{t+1}(h_{t+1}, a, j) \right\}, \quad h_t \in H_t$$

Bellman equation

$U_N(h_N) = r_N(s_N) \leftarrow$ terminal reward at $t=N$

$$u_t^*(h_t) = \sup_{\pi \in \Pi^{HD}} u_t^\pi(h_t), \quad h_t \in H_t$$

Theorem: let U_t be a solution to bellman equation then (bellman equation / optimality eqn)

$$U_t(h_t) = U_t^*(h_t) + \gamma \sum_{t+1} U_{t+1}(h_{t+1})$$

$$\text{in particular } U_t(s) = U_t^*(s) = v^*(s) = \sup_{\pi} v_{\pi}^{\pi}(s)$$

- proof:
- 1) $U_n(h_n) \geq U_n^*(h_n) + \gamma U_{n+1}(h_{n+1}) \quad \forall n \in N$
 - 2) $\forall \epsilon > 0, \exists \pi \in \Pi^{HD} \text{ s.t. } U_n^\pi(h_n) + \epsilon(N-n) \geq U_n(h_n)$
- we have to show this two to prove theorem

from 1, 2 we get: $U_n^*(h_n) + \epsilon(N-n) \geq U_n^\pi(h_n) + \epsilon(N-n) \geq U_n(h_n) \geq U_n^*(h_n)$
as $\epsilon \rightarrow 0$
 $\Rightarrow U_n^*(h_n) = U_n(h_n)$ from sandwich theorem

now let's show 1:

$$\text{at } t=N, U_N(h_N) = r_N(s_N) = U_N^\pi(h_N) \quad \forall \pi$$

$$U_N(h_N) = U_N^*(h_N) + \gamma U_{N+1}(h_{N+1})$$

if 1 holds for $n+1, \dots, N$ true for n :

$$U_n(h_n) = \sup_a \left\{ r_n(s_n, a) + \sum_{j \in S} p_n(j|s_n, a) \times U_{n+1}(h_{n+1}, a, j) \right\}$$

$$\geq \sup_a \left\{ r_n(s_n, a) + \sum_{j \in S} p_n(j|s_n, a) \times U_{n+1}^*(h_{n+1}, a, j) \right\}$$

$$\geq \sup_a \left\{ r_n(s_n, a) + \sum_{j \in S} p_n(j|s_n, a) \times U_{n+1}^\pi(h_{n+1}, a, j) \right\} + \pi \in \Pi^{HD}$$

$$\geq \sum_a \frac{q(a)}{d_n(h_n)} (r_n(s_n, a) + \sum_{j \in S} p_n(j|s_n, a) \times U_{n+1}^\pi(h_{n+1}, a, j))$$

$$= U_n^\pi(h_n) + \pi \in \Pi^{HD}$$

$\forall \pi \in \Pi^{HD}$

$$\Rightarrow U_n(h_n) \geq U_n^\pi(h_n) + \pi \in \Pi^{HD} \text{ from defn}$$

$$\Rightarrow \sup_{\pi} U_n(h_n) \geq \sup_{\pi} (U_n^\pi(h_n))$$

$$\Rightarrow U_n(h_n) \geq U_n^*(h_n)$$

now let's show 2:

for $t=N$

$$U_N^\pi(h_N) = U_N(h_N), \text{ so true } \forall \epsilon > 0$$

now let $\epsilon > 0$

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

$$r_N(s_N, d_N(h_N)) + \sum_{j \in S} p_N(j|s_N, d_N(h_N)) U_{N+1}(h_{N+1}, d_N(h_N), j) + \epsilon$$

+ N, h_N , choose $d_N(h_N)$

this is from $U_t(h_t) = \sup_{\pi} \dots$
defn of sup

\uparrow
 Action = $d_N(h_N)$
 s.t. $(\dots) + \epsilon > (\dots)$

By defn sup over action

so, given that we can find action

now this is true for N as

$$U_N^\pi(h_N) = r_N(s_N) = U_N(h_N)$$

if true for $n+1, \dots, N$ true for n :

$$\begin{aligned}
 U_p^{\pi}(h_n) &= r_n(s_n, d_n(h_n)) + \sum_{j \in S} p_n(j | s_n, d_n(h_n)) U_{n+1}^{\pi}(h_n, d_n(h_n), j) \\
 &\geq r_n(s_n, d_n(h_n)) + \sum_{j \in S} p_n(j | s_n, d_n(h_n)) U_{n+1}(h_n, d_n(h_n), j) \\
 &\geq U_n(h_n) - \varepsilon - \varepsilon(N-n+1)
 \end{aligned}$$

from defn of π we took

$$\Rightarrow U_n^{\pi}(h_n) = U_n(h_n) - \varepsilon[x + N - n - k]$$

$\forall \varepsilon > 0, \exists \pi \in \Pi^{MD}$ s.t above is true

Note: If "sup" are attained in $U_t(h_t) = \sup_{a \in A_{st}} \{r_t(s_t, a) + \sum_{j \in S} p_t(j | s_t, a) U_{t+1}(h_t, a, j)\}$ then there is an MD optimal policy

Theorem: Let U_t^* be solution to Bellman equation, then:

1) $\forall t, U_t^*(h_t)$ depends on h_t only through s_t

2) $\forall \varepsilon > 0, \exists \pi \in \Pi^{MD}$ that is ε optimal

3) If all sup are attained, then \exists an optimal MD policy

Proof: At $N, U_N^*(h_N) = r_N(s_N) \quad \forall h_N \in H_N$

so 1 holds for N
now assume it holds for $n+1, \dots, N$

$$U_n^*(h_n) = \sup_a \left\{ r_n(s_n, a) + \sum_j p_n(j | s_n, a) U_{n+1}^*(j) \right\}$$

so $U_n^*(h_n)$ depends just on s_n , 2,3 follows trivially from previous theorem and backward induction

Note: From previous 2 theorems, we get:

$$\sup_{\pi \in \Pi^{MD}} V_N^{\pi}(s) = \sup_{\pi \in \Pi^{MD}} V_N^{\pi}(s), \quad s \in S$$

Dynamic programming principle / Backward induction algo:

We are going to assume that all sup are attained

Step 1: Set $t=N$ and $U_N^*(s_N) = r_N(s_N) \quad \forall s_N \in S$

Step 2: Substitute $t-1$ in place of t and compute

$U_t^*(s)$ from bellman equation
and $a_{s,t}^* \in \arg\max_{a \in A_{st}} (\text{Bellman eqn})$

Step 3: Stop at $t=1$ and so:

$d_t(s_t) = a_{s,t}^*, s \in S, t < N$ is optimal MD policy

Note: $K^2 L^{\sum_{t=1}^N |A_t|}$ = total multiplication in above algorithm

$(L^K)^{N-1}$ = total number of MD policies $d_i: S \rightarrow A$

e.g: Sell a stock, if a stock price crosses time we sell as one of the policies
so problem is when to stop the policy

Optimal stopping problems:

$T = \{1, 2, \dots, N\}$, $N < \infty$, at everytime we can either continue or stop
 $S = S' \cup \{\Delta\}$

$$A_S = \begin{cases} \{c\} & s = \Delta \\ \{c, Q\} & s \in S' \end{cases}$$

absorption state

S' = state space of DTM C

$$\gamma_t(s, a) = \begin{cases} -f_t(s) & ; s \in S', a = c \\ g_t(s) & ; s \in S', a = Q \\ 0 & ; s = \Delta, a = c \end{cases}$$

$\gamma_N(s) = h(s) \rightarrow$ terminal reward

$$\begin{aligned} P_t(j|s, a) = & \begin{cases} p_t(j|a) & ; s \in S', a = c \\ 1 & ; s \in S', a = Q, j = \Delta \\ 1 & ; s = \Delta, j = \Delta, a = c \\ 0 & ; \text{otherwise} \end{cases} \\ t < N \end{aligned}$$

19th Aug:

optimal stopping problem:

$$S = S' \cup \{\Delta\}$$

$$A_S = \{C, Q\} \quad S \in S'$$

$$A_\Delta = \{C\} \text{ for } \Delta \text{ state}$$

$$f_t(s, a) = \begin{cases} -f_t(s) & ; s \in S', a = C \\ g_t(s) & ; s \in S', a = Q \\ 0 & ; s = \Delta, a = C \end{cases}$$

↑ price of continuing
↑ reward of quitting

$$g_t(j|s, a) = \begin{cases} p_t(j|s) & ; j \in S', s \in S', a = C \\ 1 & ; j = \Delta, s \in S', a = Q \\ 0 & ; \text{otherwise} \end{cases}$$

The secretary problem:

N candidates for job offer, $\{1, 2, \dots, N\}$, company decides to choose the candidate or reject

Objective is to maximize the prob of choosing the best candidate

$$T = \{1, 2, \dots, N\} \quad S' = \{0, 1\}$$

$$A_S = \{C, Q\}$$

↑ not the best candidate compared to last 1
↑ give offer
None person

$$f_t(s) = 0 \quad (\text{continuing No cont})$$

$$g_t(0) = 0$$

$$g_t(1) = \text{Prob}\{\text{Best candidate is in first } t\}$$

= Number of sets of size t which include the best
No. of sets of subsets E

$$\frac{n-1}{n} \binom{t-1}{t-1} = \frac{(n-1)!}{(n-t)! (t-1)!} = \frac{t}{n}$$

$$\frac{1}{\binom{n}{t}}$$

$$\Rightarrow g_t(1) = \frac{t}{N}$$

the terminal reward: $h(1) = 1$
 $h(0) = 0$

$$P_t(1|0) = \frac{1}{t+1}$$

uniformly picking one guy from $t+1$

$$P_t(1|1) = \frac{1}{t+1}$$

so $P_t(1, s) = \frac{1}{t+1}$ for $s \in \{0, 1\}$

$$P_t(0|S) = \frac{t}{1+t} \rightarrow \text{conditional prob for } S \in \{0,1\}$$

other probability same as optimal stopping problem

optimal solution for secretary problem:

let $U_t^*(1)$ = maximum probability of choosing the best candidate given that you are at the best so far (at time t)

$U_t^*(0)$ = maximum probability of choosing the best candidate given we are not at best so far (at time t)

terminal $t=N$:

$$U_N^*(1) = u(1) = 1$$

$$U_N^*(0) = h(0) = 0$$

$$U_N^*(\Delta) = 0$$

Value of choosing Q or C

$$\text{for } t < N: U_t^*(0) = \max_{Q,C} \{V^Q, V^C\}$$

$$\begin{aligned} &= \max_{Q,C} \{g_t(0) + U_{t+1}^*(\Delta), -f_t(0) + P_t(1|0)U_{t+1}^*(1) \\ &\quad + P_t(0|0)U_{t+1}^*(0)\} \\ &\text{cost to go further} \\ &\text{if we quit/continue} \\ &\text{given } s_t = 0 \\ &= \max \left\{ 0, \frac{1}{t+1} U_{t+1}^*(1) + \frac{t}{t+1} U_{t+1}^*(0) \right\} = \frac{1}{t+1} U_{t+1}^*(1) + \frac{t}{t+1} U_{t+1}^*(0) \end{aligned}$$

$$\begin{aligned} U_t^*(1) &= \max \{g_t(1) + U_{t+1}^*(\Delta), -f_t(1) + P_t(1|1)U_{t+1}^*(1) \\ &\quad + P_t(0|1)U_{t+1}^*(0)\} \\ &\text{cost to go further} \end{aligned}$$

$$\begin{aligned} &\text{if we quit/cont} \\ &\text{given } s_t = 1 \\ &= \max \left\{ \frac{t}{N}, \frac{1}{t+1} U_{t+1}^*(1) + \frac{t}{t+1} U_{t+1}^*(0) \right\} \\ &= \max \left\{ \frac{t}{N}, U_t^*(0) \right\} \end{aligned}$$

$$U_t^*(\Delta) = 0 + U_{t+1}^*(\Delta) \Rightarrow U_t^*(\Delta) = 0 \quad \forall t \in \{1, \dots, N\}$$

Prob of reaching Δ ≈ 1

so in case of $s_t = 0$, we always continue as it is maximising $U_t^*(0)$
for case of $s_t = 1$, we have 2 choices of continuing or quitting

for some $\tau \in \{1, 2, \dots, N\}$ if we integrate $t < \tau$ we continue always
after that for smallest $t > \tau$, s.t. $s_t = 1$, we do action Q
formally, suppose it is optimal to
some $\tau \in \{1, 2, \dots, N\}$ let the optimal be C

$$U_\tau^*(1) > \tau/N \quad (\text{as we are continuing we pick } U_\tau^*(1) = U_\tau^*(0))$$

$$\Rightarrow U_\tau^*(0) = U_\tau^*(1)$$

then for $t = \tau-1$

$$U_{\tau-1}^*(0) = \frac{1}{\tau} U_\tau^*(1) + \frac{\tau-1}{\tau} U_\tau^*(0)$$

$$= U_\tau^*(1) > \tau/N > \tau-1/N$$

$$\Rightarrow U_{\tau-1}^*(0) = U_{\tau-1}^*(1) > (\tau-1)/N \text{ so we continue to continue}$$

this is just to show
that if τ still true
all values less

80, optimal policy: $\exists \gamma \in \{1, 2, \dots, N\}$ s.t. when $t \leq T$, the optimal action is γ

when $t > \gamma$, if $s_t = 0 \rightarrow$ continue
 $s_t = 1 \rightarrow$ quit

Lemma: If $N > 2$, then $\gamma \geq 1$

Proof: If not true, $\gamma < 1$, or $\forall t, U_t^*(1) = \frac{t}{N} \leq U_t^*(0)$ we quit one we see $s_t = 1$

$$U_t^*(0) = \frac{1}{t+1} + \frac{t}{t+1} U_{t+1}^*(0)$$

$$\Rightarrow U_t^*(0) = \frac{1}{N} + \frac{t}{t+1} U_{t+1}^*(0)$$

$$U_N^*(0) = 0, \quad U_{N-1}^*(0) = \frac{1}{N}$$

$$U_{N-2}^*(0) = \frac{1}{N} + \frac{N-2}{N-1} \left(\frac{1}{N} \right) = \frac{N-2}{N} \left[\frac{1}{N-1} + \frac{1}{N-2} \right]$$

⋮

$$U_1^*(0) = \frac{t}{N} \left[\frac{1}{t} + \frac{1}{t+1} + \dots + \frac{1}{N-1} \right]$$

$$U_1^*(0) = \frac{1}{N} \left[1 + \frac{1}{2} + \dots + \frac{1}{N-1} \right] > \frac{1}{N} = U_1^*(1) \gg U_1^*(0) (\because \text{max value})$$

$\Rightarrow V_1^*(0) > U_1^*(0)$ this is a contradiction and so

$\gamma \geq 1$

Note: So for $N > 2$, $\gamma \geq 1$ i.e. $U_1^*(0) = U_2^*(0) = \dots = U_\gamma^*(0) \rightarrow$ continuing $t < \gamma$

$$\text{and } U_1^*(0) = U_1^*(1)$$

$$U_2^*(0) = U_2^*(1)$$

$$\text{for } t > \gamma: U_t^*(1) = \frac{t}{N}$$

$$t > \gamma: U_t^*(0) = \frac{t}{N} \left[\frac{1}{t} + \frac{1}{t+1} + \dots + \frac{1}{N-1} \right]$$

$$\cancel{\frac{t}{N}} \left[\frac{1}{t} + \dots + \frac{1}{N-1} \right] > \cancel{\frac{t}{N}}$$

$\Rightarrow \frac{1}{t} + \dots + \frac{1}{N-1} > 1$ for all $t > \gamma$ we continue no matter what

for $\frac{1}{t} + \dots + \frac{1}{N-1} < 1$ we quit if $s_t = 1$

Note: $\max \left\{ t \in \{1, \dots, N\} \mid \left[\frac{1}{t} + \frac{1}{t+1} + \dots + \frac{1}{N-1} \right] > 1 \right\} = \gamma$

$$1 \approx \left[\frac{1}{t} + \dots + \frac{1}{N-1} \right] = \int_{\gamma(N)}^{N-1} \frac{1}{x} dx = 1$$

$$\Rightarrow \log(N-1) - \log(\gamma(N)) = 1$$

$$\Rightarrow \log(N-1) - \log(e) = \log(\tau(N))$$

$$\Rightarrow \frac{N-1}{e} = \tau(N)$$

$$\text{as } N \rightarrow \infty \lim_{N \rightarrow \infty} \frac{\tau(N)}{N} = e^{-1}$$

for N large it will be $\frac{N}{e}$

22nd Aug:

Infinite horizon:

$$0 < \lambda < 1, \lim_{N \rightarrow \infty} \mathbb{E}_s^{\pi} \left[\sum_{t=1}^N \lambda^{t-1} r_t(x_t, y_t) \right] \xrightarrow{\text{Discounted reward case}} \text{we want to maximize this}$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_s^{\pi} \left[\sum_{t=1}^N r_t(x_t, y_t) \right] \xrightarrow{\text{Avg reward case}} \text{we want to maximize this}$$

Discounted case:

Expected total discounted reward = $\lim_{N \rightarrow \infty} \mathbb{E}_s^{\pi} \left[\sum_{t=1}^N \lambda^{t-1} r_t(x_t, y_t) \right]$

\hookrightarrow first assumption

$0 < \lambda < 1$ is discount factor, we assume S, A are discrete (at most countable)

time homogeneous rewards & t.p. functions i.e. $r(s, a) = r_t(s, a)$ \hookrightarrow third assumption

$(r_t(s_t, a_t))$ (transition prob functions) $P(j|s, a) = P_t(j|s, a)$ \hookrightarrow fourth assumption

$\sup_{s \in S} \sup_{a \in A_s} |r(s, a)| < \infty \hookrightarrow$ final assumption
limit always exist

$$v_{\lambda}^{\pi}(s) = \mathbb{E}_s^{\pi} \left[\sum_{t=1}^{\infty} \lambda^{t-1} r(x_t, y_t) \right]$$

we get λ as discount factor as if horizon $\sim \gamma \sim \text{geom}(1-\lambda)$ (motivation)

$$\gamma = P(\gamma = k) = \lambda^{k-1} (1-\lambda) \quad k \geq 1 \quad \text{on way } \lambda$$

$$\text{so, } \mathbb{E}_s^{\pi} \left[\mathbb{E}_{\gamma} \left(\sum_{t=1}^{\infty} r(x_t, y_t) \right) \right]$$

$$= \mathbb{E}_s^{\pi} \left[\sum_{n=1}^{\infty} \lambda^{n-1} (1-\lambda) \sum_{t=1}^n r(x_t, y_t) \right]$$

$$= \mathbb{E}_s^{\pi} \left[\sum_{t=0}^{\infty} \sum_{n=t}^{\infty} \lambda^{n-t} (1-\lambda) r(x_t, y_t) \right]$$

$$= \mathbb{E}_s^{\pi} \left[\sum_{t=0}^{\infty} \lambda^{t-1} r(x_t, y_t) \right]$$

t telescopic

Defn: (optimal policy in case of $\lambda, N \rightarrow \infty$) π^* is suitable optimal if

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

Defn: $v_{\lambda}^*(s)$ is value of MDP starting from s

$$\text{as } v_{\lambda}^{\pi}(s) = \mathbb{E}_s^{\pi} \left[\sum_{t=1}^{\infty} \lambda^{t-1} r(x_t, y_t) \right]$$

$$= \sum_{t=1}^{\infty} \lambda^{t-1} \underbrace{\mathbb{E}_s^{\pi} [r(x_t, y_t)]}_{\text{only depends on prob dist of } x_t, y_t}$$

\mathbb{P}^{π} \Rightarrow joint dist of all x_t, y_t
if we can write (x'_t, y'_t) s.t

$(x_t, y_t) \stackrel{d}{=} (x'_t, y'_t)$ then $v_{\lambda}^{\pi}(s)$ will not change even if (x'_t, y'_t) MR

Theorem: let $\pi = (d_1, d_2, \dots) \in \Pi^{HR}$, for every $s \in S, \exists \pi' = (d'_1, d'_2, \dots) \in \Pi^{MR}$ s.t

$$p^{\pi'}(X_n=j, Y_n=a | X_1=s) = p^{\pi}(X_n=j, Y_n=a | X_1=s)$$

consequently $v_{\lambda}^{\pi}(s) = v_{\lambda}^{\pi'}(s)$

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

$$q_{d_t'(j)}(a) = \sum_{\substack{\pi \in P \\ (x_1, \dots, x_{t-1}) \in H_{t-1}}} \pi(Y_t=j | X_t=j, (x_1, \dots, x_{t-1}) \in H_{t-1}, x_i=s) \quad (\text{we use induction})$$

$$\Rightarrow q_{d_t'(j)}(a) = P^\pi(Y_t=j | X_t=j, x_i=s) \quad (\text{this is our defn})$$

then eqn true for $n=1$ (trivial), if true till $n-1$ true for n :

$$\begin{aligned} P^{\pi'}(X_n=j, Y_n=a | X_1=s) &= P^{\pi'}(X_n=j | X_1=s) P^{\pi'}(Y_n=a | X_n=j, X_1=s) \\ &= P^{\pi'}(X_n=j | X_1=s) P^{\pi'}(\dots | \dots) \quad \leftarrow \text{from defn} \\ &\stackrel{?}{=} P^{\pi}(X_n=j | X_1=s) P^{\pi}(\dots | \dots) \\ \text{so if we know } P^{\pi}(X_n=j | X_1=s) &= P^{\pi}(X_n=j | X_1=s) \quad \text{then we are done} \end{aligned}$$

$$\begin{aligned} P^{\pi}(X_n=j | X_1=s) &= \sum_{i,a} P^{\pi}(X_n=j, X_{n-1}=i, Y_{n-1}=a | X_1=s) \\ &= \sum_{i,a} P^{\pi}(X_{n-1}=i, Y_{n-1}=a | X_1=s) P^{\pi}(X_n=j | X_{n-1}=i, Y_{n-1}=a, X_1=s) \\ &= \sum_a P^{\pi'}(\dots | \dots) \quad \underbrace{P^{\pi'}(\dots)}_{\substack{P(j|i,a) \text{ same} \\ \text{as } P_T(j|i,a) = P(j|i,a)}} \\ &\text{from induction} \\ &= P^{\pi'}(X_n=j | X_1=s) \end{aligned}$$

so π' depends on s

Ex: If $X \in \mathcal{D}$, then $\exists \pi'$ (depending on \mathcal{D}) s.t condition of theorem holds
 \rightarrow down down

Note: $\sup_{\pi \in \Pi_{HR}} V_\lambda^\pi(s) = \sup_{\pi \in \Pi_{MR}} V_\lambda^\pi(s)$

Policy evaluation:

We want to compute $V_\lambda^\pi(s)$ $|S|=d$

$$V_\lambda^\pi = (V_\lambda^\pi(s), s \in \mathcal{S}) \quad \text{vector of dim } d$$

$$\underset{\substack{d \in \mathcal{D}^{MD} \\ \text{decision}}} \pi(d(s)) = \pi_d(s)$$

$$\gamma_d = (r(s, d(s)), s \in \mathcal{S})$$

↑ Expected value of reward

$$r_d = \left(\sum_a r(s, a) q_{d(s)}(a), s \in \mathcal{S} \right) \quad d \in \mathcal{D}^{MD}$$

$$(P_d)_{s,j} = P(j|s, d(s)) \quad \text{matrix } d \times d \text{ as } |\mathcal{S}|=d$$

$$\|\varphi\| = \sup_{s \in S} |\varphi(s)| \quad (\text{our defn})$$

$V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ and let $\mathcal{V} : \text{set of all bounded fn from } \mathcal{S} \rightarrow \mathbb{R}$
 norm for \mathcal{V} is the sup norm

now $(V, \|\cdot\|)$ is Banach space and so $(v_\lambda^\pi \in V \text{ as } v_\lambda^\pi: f \rightarrow \mathbb{R})$
 $P: V \rightarrow V$ is a linear map
 \uparrow
 $\|P\| = \sup \left\{ \|Pv\|, \|v\| \leq 1 \right\}$

then set $BL(V) = \text{space with matrix norm s.t all maps are bounded i.e } \sup < \infty$

both V & $BL(V)$ are vector spaces and both are Banach spaces

now, $\pi = (d_1, \dots)$

$$\begin{aligned} v_\lambda^\pi(s) &= E_s^\pi \left[\sum_{t=1}^{\infty} \lambda^{t-1} \sigma(X_t, Y_t) \right] \\ &= \underbrace{\sum_a \sigma(s, a) q_{d_1(s)}(a)}_{\sigma_{d_1}(s)} + \underbrace{\lambda E_s^\pi [\sigma(X_2, Y_2)]}_{\lambda \sum_j \sigma(j, d_2(j)) P(j|s, d_1(s))} + \dots \\ &\quad = \lambda \sum_j r_{d_2(j)}(P_{d_1})_{s,j} \xleftarrow[\text{multiplication}]{\text{right}} \\ \text{the third term: } & \lambda^2 (P_{d_1} P_{d_2} r_{d_3})_s \end{aligned}$$

$$\begin{aligned} \text{so, } v_\lambda^\pi &= \sigma_{d_1} + \lambda P_{d_1} \sigma_{d_2} + \lambda^2 P_{d_1} P_{d_2} \sigma_{d_3} + \dots \\ &= \sigma_{d_1} + \lambda P_{d_1} [\sigma_{d_2} + \lambda P_{d_2} \sigma_{d_3} + \dots] \end{aligned}$$

$$\Rightarrow v_\lambda^\pi = \sigma_{d_1} + \lambda P_{d_1} \cdot v_\lambda^{\pi'}$$

where $\pi' = (d_2, d_3, \dots)$

Note: If $\pi \in \pi^{SR}$ i.e. $\pi = (d, d, \dots)$

$$\Rightarrow v_\lambda^\pi = \sigma_d + \lambda P_d v_\lambda^\pi$$

Ex: If $X_1 = y$, then $\exists \pi'$ (depending on y) s.t condition of theorem holds

Ans: Now as $X_1 \sim V$ i.e it is a random variable instead of $X_1 = s$, we make it an E value

so as $v_\lambda^\pi(s)$ known for $X_1 = s$

$\Rightarrow E(v_\lambda^\pi(X_1))$ is what we want to maximize, in all cases where we see s , so theorem still will hold

26th Aug:

$$V_\lambda^{\pi}(s) = \mathbb{E}_s^\pi \left[\sum_{t=1}^{\infty} \gamma (x_t, y_t) \right]$$

only depends on marginal so $V_\lambda^{\pi}(s) = V_\lambda^{\pi'}(s)$
for some $\pi' \in \Pi^{MR}$
(proved previously)

$$\pi = (d, d, \dots) \text{ then } V_\lambda^{\pi} = r_d + \lambda P_d V_\lambda^{\pi} \quad (\pi \in \Pi^{SR})$$

(vector wise)

Note: let L_d be a map s.t.

$$L_d: \mathcal{V} \rightarrow \mathcal{V} \quad (V_\lambda^{\pi} \in \mathcal{V})$$

$$\text{s.t. } L_d v = r_d + \lambda P_d v$$

$$r_d(s) = \sum a q_{d(s)}(a) \text{ as } \sup_{s \in S} |r_d(s)| \leq M \text{ because of bounded reward}$$

$$\Rightarrow r_d \in \mathcal{V} \text{ (bounded reward assumption)}$$

now, $\|P_d v\| \leq \|P_d\| \cdot \|v\| \quad (\because \|P\| = \sup \{ \|Pv\| : \|v\| \leq 1 \} \text{ from our defn})$

and $\|P_d\| = \sup_{S \in S} \sum_j |P(j|s)| = \frac{\|P\| \|v\|}{\|v\|}$ to make $\|v\| \leq 1$ i.e. $\frac{v}{\|v\|}$

$$\Rightarrow \|P_d v\| \leq \|v\|$$

$$\Rightarrow P_d v \in \mathcal{V}$$

as $\|v\|$ is bounded

so $P_d v$ is also bounded

thus

$$\|P\| = \sup \sum_j |P(j|s)|$$

now as $r_d \in \mathcal{V}$, $P_d v \in \mathcal{V}$ we have $r_d + \lambda P_d v \in \mathcal{V}$ (as \mathcal{V} is a vector space)

$$L_d v = r_d + \lambda P_d v \quad L_d: \mathcal{V} \rightarrow \mathcal{V}$$

and now

$$L_d V_\lambda^{\pi} = V_\lambda^{\pi}$$

Note: V_λ^{π} is a fixed point of L_d

Theorem: let $d \in D^{MR}$, then V_λ^{π} is the unique fixed point of L_d , moreover

proof:

$$\text{for some } v \in \mathcal{V} \quad \text{if } L_d v = r_d + \lambda P_d v = v$$

$$\Rightarrow v = r_d + \lambda P_d v$$

$$\Rightarrow (I - \lambda P_d)v = r_d \quad \lambda < 1$$

now as P_d is a stochastic matrix $\Rightarrow \max |\mu| = 1$

(μ = eigenvalues of P_d)

as $\|P_d\| = \max(|\mu|)$ by defn of our eigenvalues and norm

$$\|P_d\| = 1$$

$$\Rightarrow 1 = \max(|\mu|)$$

and so $\max |\lambda \mu| < 1$ as $\lambda < 1$

& so 0 cannot be eigenvalue of $I - \lambda P_d$
as 0 is not an eigenvalue of $I - \lambda P_d \Rightarrow I - \lambda P_d$ is invertible
as $\det(I - \lambda P_d) = \lambda_1 \lambda_2 \dots \lambda_n \rightarrow$ all eigenvalues non-zero

$$\neq 0$$

so $I - \lambda P_d$ is invertible
and so $v = (I - \lambda P_d)^{-1} r_d$

↓
unique

Defn: (spectral radius) let $T: \mathcal{V} \rightarrow \mathcal{V}$, then $\sigma(T) = \sup \{ |\mu| \mid (\mu I - T)^{-1} \text{ is not invertible} \}$

Note: $\sigma(T) = \lim_{n \rightarrow \infty} \|T^n\|^{\frac{1}{n}}$

Note: If $\sigma(\tau) < 1$ then $I - \tau$ is invertible and $(I - \tau)^{-1} = \sum_{n=0}^{\infty} \tau^n$

now, $\sigma(\lambda P_d) = \lambda < 1$ so, $(I - \lambda P_d)$ is invertible

$(\because \sigma(\lambda P_d) = \lim_{n \rightarrow \infty} \lambda^n / \lambda^n = \lambda)$

we can compute $(I - \tau)^{-1}$ using same

alternate way to know $(I - \lambda P_d)$ is invertible

Now, $\sup_{a \in A^S} \left\{ r(s, a) + \lambda \sum_{j \in S} p(s, a, j) V(j) \right\} = V(s)$, is the bellman equation for infinite horizon problem also called optimality equation

Now, $L: V \rightarrow V$ be the bellman operator

$$L\varphi = \sup_{d \in D^{MD}} \{ (r_d + \lambda P_d \varphi) \}$$

$$\text{s.t. } L\varphi(s) = \sup_{d \in D^{MD}} \sup_{a \in A^S} ((r_d + \lambda P_d \varphi)(s))$$

Note: $L\varphi = \sup_{d \in D^{MR}} \{ r_d + \lambda P_d \varphi \}$ from previous lemma (as sup does not matter)

suppose $V \geq L\varphi$ & component

then we want to show $V \geq V_\lambda^*$

as $V \geq L\varphi \geq r_{d_1} + \lambda P_{d_1} \varphi$ from defn of L

$$\Rightarrow V \geq r_{d_1} + \lambda P_{d_1} (r_{d_2} + \lambda P_{d_2} \varphi)$$

⋮

$$\geq r_{d_1} + \lambda P_{d_1} r_{d_2} + \dots + \lambda^{n-1} P_{d_1} \dots P_{d_{n-1}} r_{d_n} + \lambda^n P_{d_1} \dots P_{d_n} \varphi$$

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

then

$$V_\lambda^\pi = r_{d_1} + \lambda P_{d_1} r_{d_2} + \dots$$

$$\forall \pi \in \Pi^{MR}$$

$$\text{then } V \geq V_\lambda^\pi - \sum_{k=n}^{\infty} \lambda^k P_{d_1} \dots P_{d_{k-1}} r_{d_{k+1}} + \lambda^n P_{d_1} \dots P_{d_n} \varphi$$

$$\Rightarrow V - V_\lambda^\pi \geq \lambda^n P_{d_1} \dots P_{d_n} \varphi - \sum_{k=n}^{\infty} \lambda^k P_{d_1} \dots P_{d_{k-1}} r_{d_{k+1}}$$

this is true $\forall n$ as $P_{d_1} \times P_{d_2} \dots$ is a stochastic matrix

now $| \lambda^n P_{d_1} \dots P_{d_n} \varphi | \leq \lambda^n \times \| \varphi \| \leq \lambda^n \| \varphi \|$ as $n \rightarrow \infty \lambda^n \rightarrow 0$ so

given $\epsilon > 0$, $\exists N$ s.t. $-\sum e \leq \lambda^n \leq \sum e$

all one vector $e = (1, 1, \dots)$

$$\sum \lambda^k P_{d_1} \dots P_{d_{k-1}} r_{d_{k+1}} \leq M \underbrace{\sum \lambda^k P_{d_1} \dots P_{d_{k-1}} e}_{e \text{ as stochastic matrix}}$$

$$= M e \frac{\lambda^n}{1-\lambda}$$

$$\exists N' > N \text{ s.t. } \frac{\lambda^n M}{1-\lambda} < \epsilon \text{ as } \lambda^n \rightarrow 0 \text{ as } n \rightarrow \infty$$

so for sufficient large n we get:

$$V - V_\lambda^\pi \geq -2\epsilon e \quad \forall \epsilon > 0$$

then as $\epsilon \rightarrow 0$

$$\Rightarrow V - V_\lambda^\pi \geq 0$$

$$\Rightarrow V \geq V_\lambda^\pi \quad \forall \pi \in \Pi^{MR}$$

$$\Rightarrow v \geq \sup_{\pi \in \Pi^M} v_\lambda^\pi = v_\lambda^*$$

so if $v > v_\lambda^*$ $\Rightarrow v > v_\lambda^*$

now, we want to find $v \leq Lv$ then does $v \leq v_\lambda^*$, for this:

let $\epsilon > 0$, then

$v \leq r_d + \lambda Pd v + \epsilon \cdot e$ for some d (from defn of supremum)

$$\Rightarrow (I - \lambda Pd)v \leq r_d + \epsilon \cdot e$$

$$\Rightarrow v \leq (I - \lambda Pd)^{-1}(r_d + \epsilon \cdot e) \quad (\because (I - \lambda Pd)^{-1} = \sum_{n=0}^{\infty} (\lambda Pd)^n \text{ all positive so positive and})$$

$$\text{as } (I - \lambda Pd)^{-1}(r_d + \epsilon \cdot e) = v(d, d, \dots) + \frac{\epsilon \cdot e}{1-\lambda} \quad (\text{first term is straight from prev theorem, last term from})$$

$$\Rightarrow v \leq \sup_{\pi \in \Pi^M} v_\lambda^\pi + \frac{\epsilon \cdot e}{1-\lambda} \quad \text{as } \epsilon \text{ is arbitrary}$$

$$\Rightarrow v \leq v_\lambda^*$$

$$\sum_{n=0}^{\infty} (\lambda^n) P_d^n \cdot e = \frac{e}{1-\lambda}$$

as P_d^n is stochastic

Theorem: (a) If $v \geq Lv$ then $v \geq v_\lambda^*$
 (b) If $v \leq Lv$ then $v \leq v_\lambda^*$
 (c) If $v = Lv$ then $v = v_\lambda^*$

Proof: The proof follows trivially from calculations done above

Note: (c) part does not imply L has a fixed point, as we know if $Lv = v$ then $v = v_\lambda^*$ (we know if L has a fixed point then it is v_λ^*)

Defn: $T: V \rightarrow V$ is said to be a contraction if $\exists \lambda \in \mathbb{R}$ s.t $\|Tu - Tv\| \leq \lambda \|u - v\|$ $\forall u, v \in V$

Theorem: (Banach fixed point theorem) Let $T: V \rightarrow V$ be a contraction then:

(a) \exists a unique v^* s.t $Tv^* = v^*$

(b) Starting with any $v^0 \in V$, $v^{n+1} = T v^n = T^2 v^{n-1} = \dots = T^{n+1} v^0$ for $n \geq 1$ converges to v^*

2nd Sept:

$$\text{we have } V_x^\pi(s) = \mathbb{E}_s^\pi \left[\sum_{t=1}^{\infty} \lambda^{t-1} r(s_t, a_t) \right]$$

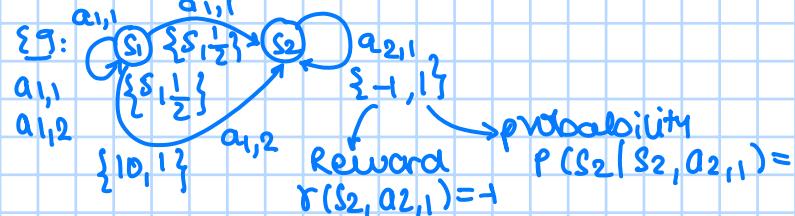
(d, d, ...)

$$L_d V = r_d + \lambda P_d V$$

& we get $V = \sup_{\pi \in DMD} \left\{ r_d + \lambda P_d V \right\}$ Bellman equation

$\overbrace{L_d V}$ from theorem done

$$\text{so if } V \text{ is a fixed point then } V = V^*$$



$$a(S_1) = a_{1,1}$$

$\pi = (d, d, \dots)$ then we want to find $V_\lambda^\pi = V$
 $d(S_2) = a_{2,1} \rightarrow \text{only one possible action}$

$$V(S_2) = -1 + \lambda \begin{pmatrix} 0 \\ 1 \end{pmatrix} (V(S_1), V(S_2))$$

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

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

$$\begin{aligned} \text{now, } V(S_1) &= s + \lambda \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} (V(S_1), V(S_2)) \\ &= 5 + \lambda \left[\frac{V(S_1)}{2} + \frac{1}{2} V(S_2) \right] \\ \Rightarrow V(S_1) &= \frac{10 - 11\lambda}{(2-\lambda)(1-\lambda)} \end{aligned}$$

$$V_x^\pi = (I - \lambda P_d)^{-1} r_d \text{ can also be used}$$

Note: In above example for optimality equation as only $a_{2,1}$
 we get $V_\lambda^*(S_2) = \frac{-1}{1-\lambda}$

$$\text{and for } S_1: V(S_1) = \max \left\{ s + \lambda \left[\frac{1}{2} V(S_1) + \frac{V(S_2)}{2} \right], 10 + \lambda V(S_2) \right\}$$

$$\text{if } \lambda = 0 \text{ then } V(S_1) = \max \{ s, 10 \} = 10$$

$$\text{so } V(S_2) = -1$$

$V(S_1) = 10$ optimal action is $a_{1,2}$

$$\lambda = \frac{1}{2}: V(S_2) = -2$$

$$V(S_1) = \max \left\{ s - \frac{1}{2} + \frac{1}{4} V(S_1), 9 \right\}$$

$$V(S_1) = 9 \text{ as } \frac{1}{4} V(S_1) = \frac{1}{2} \Rightarrow V(S_1) = 6 \text{ in } a_{1,1} \text{ case}$$

so for $\lambda = 1/2$ $a_{1,2}$ is optimal action

one: Find an optimal action at s_i as a function of $\lambda \rightarrow$ down down

now if $v = Lv$ then $v = v^*$ (theorem done previously)

Defn: $T: V \rightarrow V$ is said to be a contraction if $\exists \lambda < 1$ s.t. $\|Tu - Tv\| \leq \lambda \|u - v\|$
 $\forall u, v \in V$

Theorem: (Banach fixed point theorem) let $T: V \rightarrow V$ be a contraction then:
where V is a banach space

(a) \exists a unique v^* s.t. $Tv^* = v^*$

(b) Starting with any $v^0 \in V$, $v^{n+1} = T v^n = T^2 v^{n-1} = \dots = T^{n+1} v^0$
for $n \geq 1$ converges to v^*

$$\|v_n - v^*\| \rightarrow 0 \text{ as } n \rightarrow \infty$$

proof: we fix m, n $\xleftarrow{n} \xrightarrow{m+n}$ we want to compare

$$\|v^{n+m} - v^n\| \leq \sum_{k=0}^{m-1} \|v^{n+k+1} - v^{n+k}\| \text{ by triangle inequality}$$

$$= \sum_{k=0}^{m-1} \|T^{n+k} v^1 - T^{n+k} v^0\|$$

$$\leq \sum_{k=0}^{m-1} \lambda^{n+k} \|v^1 - v^0\| \quad (\because \text{By contraction property})$$

$$= \lambda^n \|v^1 - v^0\| \sum_{k=0}^{m-1} \lambda^k$$

$$= \lambda^n \|v^1 - v^0\| \left[\frac{(1-\lambda^m)}{(1-\lambda)} \right]$$

$$\Rightarrow \|v^{n+m} - v^n\| \leq \lambda^n \|v^1 - v^0\| \left[\frac{(1-\lambda^m)}{(1-\lambda)} \right]$$

as $\lambda < 1$, $\forall \epsilon > 0, \exists N$ s.t.

$$\|v^{n+m} - v^n\| < \epsilon \quad \forall m \geq 1, n \geq N$$

so $\{v^n\}$ is a cauchy sequence and as V is a banach space (every cauchy seqn converges)

→ this is from completeness

$$\text{so } \|v^n - v^*\| \rightarrow 0 \text{ as } n \rightarrow \infty$$

and now we want to show $Tv^* = v^*$

$$\begin{aligned} & \|Tv^* - v^*\| \\ &= \|Tv^* - v^n + v^n - v^*\| \\ &\leq \|Tv^* - v^n\| + \|v^n - v^*\| \\ &= \|Tv^* - T v^{n-1}\| + \|v^n - v^*\| \\ &= \lambda \|v^* - v^{n-1}\| + \|v^n - v^*\| \\ &\text{as } n \rightarrow \infty \quad \|Tv^* - v^*\| \rightarrow 0 \\ &\Rightarrow Tv^* = v^* \end{aligned}$$

also v^* is unique as if not true $\exists u^* \text{ s.t. } Tu^* = u^*$

$$Tv^* = v^*$$

$$\& u^* \neq v^*$$

$$\|Tu^* - Tv^*\| \leq \lambda \|u^* - v^*\| < \|u^* - v^*\|$$

$$\Rightarrow \|u^* - v^*\| < \|u^* - v^*\| \text{ this is a contradiction}$$

$$\& \text{so, } u^* = v^*$$

∴ v^* is a unique fixed point

Note: $L: \mathcal{V} \rightarrow \mathcal{V}$ has to be a contraction for us to have a fixed point

$$L\vartheta = \sup_{d \in DMD} \{r_d + \lambda p_d \vartheta\}$$

Theorem: Suppose all sup are attained, then L is a contraction. By all sup we mean for all $\vartheta \in \mathcal{V}$, sup is attained

Proof: We want to show that $\|L\vartheta - L\vartheta'\| \leq \lambda \|\vartheta - \vartheta'\|$

but is same λ as our objective function

lets fix a state s , for some $\vartheta, \vartheta' \in \mathcal{V}$, let

$$L\vartheta(s) \geq L\vartheta'(s)$$

fixed s

$$\text{let } a_s^* \in \arg\max_{a \in A_s} \{r(s, a) + \lambda \sum_{j \in S} p(j|s, a) \vartheta(j)\}$$

$$L\vartheta(s) \geq \left\{ r(s, a_s^*) + \lambda \sum_{j \in S} p(j|s, a_s^*) \vartheta(j) \right\} \quad (\because \text{plugging } a_s^* \text{ for } \vartheta(s))$$

$$\text{where } L\vartheta(s) - L\vartheta'(s) \geq 0$$

$$(L\vartheta(s) - L\vartheta'(s)) \leq \left\{ r(s, a_s^*) + \lambda \sum_{j \in S} p(j|s, a_s^*) \vartheta(j) \right\} \\ - \left\{ r(s, a_s^*) + \lambda \sum_{j \in S} p(j|s, a_s^*) \vartheta'(j) \right\}$$

$$\Rightarrow L\vartheta(s) - L\vartheta'(s) \leq \lambda \sum_{j \in S} p(j|s, a_s^*) [r(j) - \vartheta'(j)] \\ \leq \lambda \left(\sum_{j \in S} p(j|s, a_s^*) \right) \|r - \vartheta'\| \quad \begin{matrix} \text{taking max of } j \in S \\ \text{as } \mathcal{V} \text{ is equipped with sup norm} \end{matrix}$$

$$\text{so, } |L\vartheta(s) - L\vartheta'(s)| \leq \lambda \|\vartheta - \vartheta'\| \quad \forall s \in \mathcal{S}$$

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

$\Rightarrow L$ is a contraction

Ex: Show that even sup need not be attained, still L is a contraction \Rightarrow done down

Note: $L_d \vartheta = r_d + \lambda p_d \vartheta$, we can show that L_d is also a contraction and from banach fixed point theorem we get \exists unique fixed point

now, we showed that $L: \mathcal{V} \rightarrow \mathcal{V}$, then \exists a unique fixed point

$$L\vartheta^* = \vartheta^* \quad \& \quad \vartheta^* = \vartheta_\lambda^* \quad (\text{from prev theorem})$$

where $\vartheta_\lambda^* = \sup_{\pi} \vartheta_\lambda^\pi$

Theorem: A policy $\pi^* \in \Pi^{HR}$, then it is optimal iff $\vartheta_\lambda^{\pi^*}$ is a solution to the optimality equations

Proof: (\Rightarrow) As π^* is optimal, $\vartheta_\lambda^{\pi^*} = \vartheta_\lambda^* = L\vartheta_\lambda^*$ (By theorem done before)

$$\Rightarrow \vartheta_\lambda^{\pi^*} = L\vartheta_\lambda^{\pi^*}$$

$\Rightarrow \vartheta_\lambda^{\pi^*}$ is solution to optimality eqn

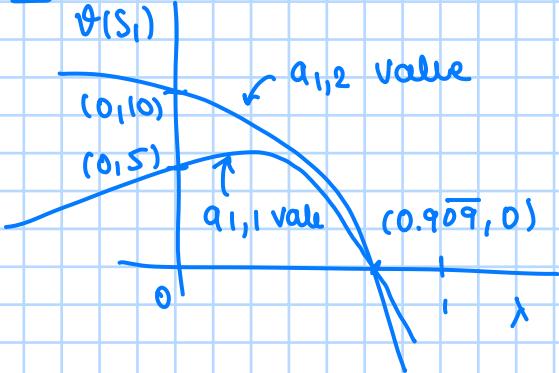
(\Leftarrow) As L has unique fixed point

$$\text{if } \vartheta_\lambda^{\pi^*} = L\vartheta_\lambda^{\pi^*}$$

$\Rightarrow \vartheta_\lambda^{\pi^*} = \vartheta_\lambda^*$ i.e. π^* is optimal as from prev theorem ϑ_λ^* is only fixed point

Ex: Find an optimal action at s_1 as a function of λ

Ans:



for $\lambda > 0.909$ $a_{1,1}$ is optimal action
or else $a_{1,2}$ is optimal

Ex: Show that max sup need not be attained, still L^\dagger is a contraction

Ans: similar to max sup is attained

for some $U, V \in \mathcal{V}$ and some $s \in S$
wlog $L(V)(s) \geq L(U)(s)$ then

and as sup need not be attained $\forall \varepsilon > 0 \exists a^* \in A$ s.t.

$$r(s, a^*) + \lambda \sum_{j \in S} P(j|s, a^*) V(j) + \varepsilon > \sup_{a \in A} \left\{ r(s, a) + \lambda \sum_{j \in S} P(j|s, a) V(j) \right\}$$

$$\text{as } L_U(s) \geq r(s, a^*) + \lambda \sum_{j \in S} P(j|s, a^*) U(j)$$

$$\Rightarrow |L_V(s) - L_U(s)| \leq r(s, a^*) + \lambda \sum_{j \in S} P(j|s, a^*) V(j) + \varepsilon - r(s, a^*) - \lambda \sum_{j \in S} P(j|s, a^*) U(j)$$

$$\Rightarrow |L_V(s) - L_U(s)| \leq \lambda \left(\sum_{j \in S} P(j|s, a^*) \right) \|V - U\| + \varepsilon$$

↓
sup norm

$$\Rightarrow |L_V(s) - L_U(s)| \leq \lambda \|V - U\| + \varepsilon$$

true for all s

$$\therefore \|L_V - L_U\| \leq \lambda \|V - U\| + \varepsilon$$

↑ sup norm and so for $\varepsilon \rightarrow 0$
we get: $\|L^\dagger - L\| \leq \lambda \|V - U\|$ so L is a contraction

9th Sept:

so far, $L: V \rightarrow V$ is a bellman operator, showed that L is a contraction and hence from banach fixed point theorem:

we cannot comment on uniqueness of policy π , only $v_\lambda^* \in V$

\exists unique fixed point $v^* \in V$ for L
s.t. $v^* = Lv^*$
 $Lv_\lambda^* = v_\lambda^*$

where $v_\lambda^* = \sup_d \{ r_d + \lambda P_d v_\lambda^* \}$
Bellman

Defn: $d^* \in D^{MD}$ is said to be covering if $d^* \in \arg\max_{d \in D^{MD}} \{ r_d + \lambda P_d v_\lambda^* \}$

where if we find, $a_s^* \in \arg\max_{a \in A_s} \{ r(s, a) + \lambda \sum_j p(j|s, a) v_\lambda^*(j) \}$
i.e $d^*(s) = a_s^*$

Theorem: Suppose all suprema are attained, then ($L: V \rightarrow V$ all sups attained)

a) \exists a covering rule d^*

b) If d^* is covering, then $\pi = (d^*, d^*, \dots) \in \Pi^{SD}$ is optimal

Proof: a. is trivial as all suprema are attained, so d^* exist

b. It suffices to show $v_\lambda^\pi = v_\lambda^*$, for this

as $v_\lambda^* = Lv_\lambda^* = \sup_d \{ r_d + \lambda P_d v_\lambda^* \}$

\leftarrow For stationary policy $= r_{d^*} + \lambda P_{d^*} v_\lambda^* = L_{d^*}(v_\lambda^*)$
 $L_d: V \rightarrow V$ s.t.

$\text{so, } L_{d^*}(v_\lambda^*) = v_\lambda^* \quad L_d(v) = r_d + \lambda P_d v$

$\Rightarrow v_\lambda^*$ is a fixed point for L_{d^*}

we have shown that L_d has a unique fixed point v_λ^π for $\pi = (d^*, d^*, \dots)$

so $v_\lambda^* = v_\lambda^\pi$ for $\pi = (d^*, d^*, \dots)$

so as $v_\lambda^\pi = v_\lambda^* \Rightarrow \pi = (d^*, d^*, \dots)$ is optimal and $\pi \in \Pi^{SD}$

Note: In particular, it means that if all sups are attained in $L: V \rightarrow V$ then:

$\sup_{\pi \in \Pi^{HR}} v_\lambda^\pi(s) = \sup_{\pi \in \Pi^{SD}} v_\lambda^\pi(s)$ (as shown $\sup_{\pi \in \Pi^{HR}} v_\lambda^\pi(s) = v_\lambda^*(s)$)

Now lets see when all sups are not attained

Theorem: For any $\epsilon > 0$, \exists an ϵ optimal SD policy

Proof: we have $Lv_\lambda^* = v_\lambda^*$ as unique fixed point of $L: V \rightarrow V$
given $\epsilon > 0$, $\exists d_\epsilon \in D^{MD}$ s.t.

\leftarrow From sup property $r_{d_\epsilon} + \lambda P_{d_\epsilon} v_\lambda^* \geq Lv_\lambda^* - \epsilon$ $\epsilon = (\cdot)$ all
 $\underbrace{\epsilon}_{\text{vector}}$ one

considering ϵ as $\epsilon(1-\lambda)$

$\Rightarrow r_{d_\epsilon} + \lambda P_{d_\epsilon} v_\lambda^* \geq v_\lambda^* - \epsilon(1-\lambda) \quad (\because Lv_\lambda^* = v_\lambda^*)$

$Lv = \sup_{d \in D^{MD}} \{ r_d + \lambda P_d v \}$

Defn of Bellman operator

$$\Rightarrow r_{d_\xi} + \varepsilon e(1-\lambda) \geq (I - \lambda P_{d_\xi}) \vartheta_\lambda^*$$

as $(I - \lambda P_{d_\xi})^\dagger$ multiplied preserves inequality (\because done before)

$$\Rightarrow \vartheta_\lambda^* \leq \underbrace{(I - \lambda P_{d_\xi})^\dagger}_{V_\lambda^\pi \text{ mean } \pi = (d_\xi, d_\xi, \dots)} r_{d_\xi} + \varepsilon(1-\lambda)(I - \lambda P_{d_\xi})^\dagger e$$

$$V_\lambda^\pi \text{ mean } \pi = (d_\xi, d_\xi, \dots) \in \Pi^{SD}$$

$$\text{and } (I - \lambda P_{d_\xi})^\dagger e = \sum_{n=0}^{\infty} (\lambda P_{d_\xi})^n e$$

$$= (1 + \lambda + \lambda^2 + \dots) e \\ = \frac{1}{1-\lambda} e$$

$$\text{so, } \vartheta_\lambda^* \leq V_\lambda^\pi + \varepsilon \cdot e$$

thus mean, $(d_\xi, d_\xi, \dots) \in \Pi^{SD}$ is an ε optimal policy

Note: $\sup_{\pi \in \Pi_{MR}} \vartheta_\lambda^\pi(s) = \sup_{\pi \in \Pi^{SD}} \vartheta_\lambda^\pi(s)$ in all cases ($\because \varepsilon \rightarrow 0$, we get this)

Value iteration algorithm:

We want to algorithmically find an optimal policy given r_d, P_d, λ
we will assume all suprema are attained

Algorithm:

1. pick $\vartheta^0 \in V$, fix $\varepsilon > 0$, set $n=0$

2. while $(\|\vartheta^n - L\vartheta^n\| > \varepsilon(1-\lambda))$:

$$\vartheta^{n+1} = L\vartheta^n$$

$n = n+1$

$$\vartheta^{n+1}(s) = \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{j \in S} P(j|s, a) \vartheta^n(j) \right\}$$

3. as $\|\vartheta^n - L\vartheta^n\| < \varepsilon \left(\frac{1-\lambda}{2\lambda} \right)$ now, $\vartheta^{n+1} = L\vartheta^n$

$$\text{choose } d_\xi \in \arg\max \left\{ r_d + \lambda P_d \vartheta^{n+1} \right\}$$

4. return (ϑ^{n+1}, d_ξ)

We now have to show ① ϑ^{n+1} is close to ϑ_λ^*

② d_ξ is ε optimal

Theorem: For the above algorithm, given $\{\vartheta^n\}$ by $\vartheta^{n+1} = L\vartheta^n$

a) $\vartheta^n \rightarrow \vartheta_\lambda^*$ as $n \rightarrow \infty$

b) $\exists N \text{ s.t. } \|\vartheta^{n+1} - \vartheta^n\| < \varepsilon \left(\frac{1-\lambda}{2\lambda} \right) \text{ holds } \forall n \geq N$

c) $\Pi^\varepsilon = (d_\xi, d_\xi, \dots) \in \Pi^{SD}$ is ε -optimal

d) $\|\vartheta^{n+1} - \vartheta_\lambda^*\| < \frac{\varepsilon}{2}$ whenever $\|\vartheta^{n+1} - \vartheta^n\| < \varepsilon \left(\frac{1-\lambda}{2\lambda} \right)$ holds

Proof: a) is trivial from Banach fixed point theorem

b) from Banach fixed point, as $\{\vartheta^n\}$ is a Cauchy sequence,
as V is a Banach space, from defn $\forall \varepsilon > 0, \exists \text{ such } N$

$$\text{for c) and d), } \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta_\lambda^*\|$$

$$= \|\vartheta_\lambda^{\pi^\varepsilon} + \vartheta^{n+1} - \vartheta^{n+1} - \vartheta_\lambda^*\|$$

$$\leq \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1}\| + \|\vartheta^{n+1} - \vartheta_\lambda^*\| \quad \text{--- ①}$$

where $\vartheta_\lambda^{\pi^\varepsilon} = L_{d_\varepsilon} \vartheta_\lambda^{\pi^\varepsilon}$ from L_{d_ε} constraining
 $L\vartheta^{n+1} = L_{d_\varepsilon} \vartheta^{n+1}$ from defn of L and L_{d_ε}
calculation of d_ε showing this $< \varepsilon/2$
will prove part d

$$\begin{aligned} \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1}\| &= \|L_{d_\varepsilon} \vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1} + L\vartheta^{n+1} - L\vartheta^{n+1}\| \\ &\leq \|L_{d_\varepsilon} \vartheta_\lambda^{\pi^\varepsilon} - \underbrace{L\vartheta^{n+1}}_{L_{d_\varepsilon} \vartheta^{n+1}}\| + \underbrace{\|L\vartheta^{n+1} - L\vartheta^n\|}_{< \lambda \|\vartheta^{n+1} - \vartheta^n\|} \\ &\leq \lambda \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1}\| + \lambda \|\vartheta^{n+1} - \vartheta^n\| \end{aligned}$$

from contraction
property of L_{d_ε}

from contraction
property of L

$$\Rightarrow (1-\lambda) \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1}\| < \lambda(\varepsilon) \left(\frac{1-\lambda}{2\lambda} \right) \quad (\because \text{from algorithm } \|\vartheta^{n+1} - \vartheta^n\| < \varepsilon \left(\frac{1-\lambda}{2\lambda} \right))$$

$$\Rightarrow \|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta^{n+1}\| < \frac{\varepsilon}{2} \quad \text{--- ②}$$

so, now from ②, ① becomes: $\|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta_\lambda^*\| < \frac{\varepsilon}{2} + \|\vartheta^{n+1} - \vartheta_\lambda^*\|$

$$\text{now, } \|\vartheta^{n+1} - \vartheta_\lambda^*\| = \|\vartheta^{n+1} - L\vartheta_\lambda^*\| \quad (\because L\vartheta_\lambda^* = \vartheta_\lambda^*)$$

$$\begin{aligned} &\leq \underbrace{\|\vartheta^{n+1} - L\vartheta^n\|}_{L\vartheta^n} + \|\underbrace{L\vartheta^n - L\vartheta^{n+1}}_{L\vartheta^{n+1}}\| + \|\underbrace{L\vartheta^{n+1} - L\vartheta_\lambda^*}_{L\vartheta_\lambda^*}\| \\ &\leq \lambda \|\vartheta^n - \vartheta^{n+1}\| + \lambda \|\vartheta^{n+1} - \vartheta_\lambda^*\| \end{aligned}$$

$$\begin{aligned} \Rightarrow \|\vartheta^{n+1} - \vartheta_\lambda^*\| &\leq \frac{\lambda}{1-\lambda} \|\vartheta^n - \vartheta^{n+1}\| \\ &< \frac{\varepsilon}{2} \quad (\because \text{stopping condition}) \quad \text{--- ③} \end{aligned}$$

so from ②, ③, ① becomes:

$$\|\vartheta_\lambda^{\pi^\varepsilon} - \vartheta_\lambda^*\| < \varepsilon \quad \text{and} \quad \|\vartheta^{n+1} - \vartheta_\lambda^*\| < \frac{\varepsilon}{2}$$

12th Sept:

Value Iteration:

$$v^{n+1} = L v^n$$

$$\|v^{n+1} - v^n\| \leq \frac{\lambda}{1-\lambda} (1-\lambda)$$

$$d \in \arg \max_d \{r_d + \lambda p_d v^{n+1}\}$$

Theorem: (a) If $u, v \in V$, $v \geq u$ (for all components) then $Lv \geq Lu$

(b) If $Lv^N \geq v^N$ for some N , then $v^{N+m+1} \geq v^{N+m} \forall m \geq 0$
and similarly

$Lv^N \leq v^N$ for some N , then $v^{N+m+1} \leq v^{N+m} \forall m \geq 0$

Proof: (a) Let $d \in \arg \max_d \{r_d + \lambda p_d u\}$, then

$$Lu = r_d + \lambda p_d u$$

and as p_d has non-negative entries (transition matrix)
 $p_d u \leq p_d v$

$$\Rightarrow Lu = r_d + \lambda p_d u \leq r_d + \lambda p_d v \leq \max_d \{r_d + \lambda p_d v\} = Lv$$

$$\Rightarrow Lu \leq Lv$$

$$(b) v^{N+m+1} = L^{m+1} v^N$$

$$= L^m (Lv^N)$$

as $Lv^N \geq v^N$ for some N

$$\Rightarrow L(Lv^N) \geq Lv^N \quad (\text{by part (a)})$$

:

$$\Rightarrow L^m (Lv^N) \geq Lv^N$$

$$\Rightarrow v^{N+m+1} \geq v^{N+m} \forall m \geq 0$$

Ex: Show the \leq part of (b) \rightarrow done done

let $r(s,a) \geq 0$, $Lv = \max_d \{r_d + \lambda p_d v\}$, then we would make $v^0 = 0$ as

$$Lv^0 = \max_d \{r_d + 0\} \geq 0$$

$$\Rightarrow Lv^0 \geq v^0$$

$$\Rightarrow v^1 \geq v^0$$

so, $v^0 \uparrow v^* \uparrow v^N$ is monotonically inc to v^*

similarly if $r_d \leq 0$ then $v^0 = 0$ gives
 $v^N \downarrow v^* \downarrow v^N$ is monotonically dec to v^*

Rate of convergence:

$\{v^n\}$ is a seq of vectors, then $v^n \rightarrow v^*$ if $\|v^n - v^*\| \rightarrow 0$ as $n \rightarrow \infty$

Defn: $\{v^n\}$ long at order (at least α) if $\alpha > 0$ and $\|v^{n+1} - v^*\| \leq K \|v^n - v^*\|^{\alpha} \forall n$
for some constant $K > 0$

↑ globally

Defn: (rate of conv) smallest sum K is called rate of convergence

Note: $\alpha = 1$: linear convergence (order 1)

Defn: (Asymptotic rate of convergence) $\lim_{n \rightarrow \infty} \left(\frac{\|\vartheta^n - \vartheta^*\|}{\|\vartheta^0 - \vartheta^*\|} \right)^{\frac{1}{n}}$ given order = 1
we call it AARC

$\{\vartheta^n\}$ converges $O(f(n))$ if $\lim_{n \rightarrow \infty} \frac{\|\vartheta^n - \vartheta^*\|}{f(n)} < \infty$

and we say $\vartheta^n = \vartheta^* + O(f(n))$

Note: If $f(n) = \beta^n$, $\beta < 1$, then we say geometric convergence at ratio β

Theorem: $\{\vartheta^n\}$ from value iteration, then the following global conv. property holds:

- (i) conv. is linear with rate λ
- (ii) AARC $\approx \lambda$
- (iii) convergence is $O(\lambda^n)$

Proof:

$$(i) \|\vartheta^{n+1} - \vartheta_\lambda^*\| = \|L\vartheta^n - L\vartheta_\lambda^*\|$$

$$\leq \lambda \|\vartheta^n - \vartheta_\lambda^*\|, \forall n$$

so, conv is atleast linear, if $\vartheta^0 = \vartheta_\lambda^* + K \cdot e$, $K \neq 0$ forall one vector

$$\text{then } \vartheta^1 = \max_d \{r_d + \lambda p_d (\vartheta_\lambda^* + K \cdot e)\}$$

$$= \max_d \{r_d + \lambda p_d \vartheta_\lambda^* + \lambda K \cdot e\}$$

$$= \max_d \{r_d + \lambda p_d \vartheta_\lambda^*\} + \lambda K \cdot e$$

$$\vartheta^1 = \lambda K \cdot e + \vartheta_\lambda^*$$

$$\Rightarrow \vartheta^1 - \vartheta_\lambda^* = \lambda(\vartheta^0 - \vartheta_\lambda^*) \quad \forall n \text{ with above choice of } \vartheta^0$$

so, $\alpha = 1$ and rate = λ

$$(ii) \|\vartheta^n - \vartheta_\lambda^*\| \leq \lambda^n \|\vartheta^0 - \vartheta_\lambda^*\| \quad \text{--- (1)}$$

$$\text{so, } \overline{\lim}_{n \rightarrow \infty} \left(\frac{\|\vartheta^n - \vartheta_\lambda^*\|}{\|\vartheta^0 - \vartheta_\lambda^*\|} \right)^{\frac{1}{n}} \leq \lambda \quad \text{but from (a) } \vartheta^0 = \vartheta_\lambda^* + K \cdot e \quad K \neq 0$$

makes AARC $\approx \lambda$

$\underbrace{\text{AARC}}$

so, AARC = λ

$$(iii) \overline{\lim}_{n \rightarrow \infty} \frac{\|\vartheta^n - \vartheta^*\|}{\lambda^n} = \|\vartheta^0 - \vartheta^*\| < \infty \quad \text{so, converges } O(\lambda^n)$$

gauss-seidel value iteration:

$\{s_1, \dots, s_N\}$ states

algorithm:

① Fix $\epsilon > 0$, $\vartheta^0 \in V$, $n > 0$

② $\vartheta^{n+1}(s_j) = \max_{a \in A_{s_j}} \{r(s_j, a) + \lambda \left[\sum_{i < j} p(s_i | s_j, a) \vartheta^{n+1}(i) + \sum_{i > j} p(s_i | s_j, a) \vartheta^n(i) \right]\}$

sequentially for $j = 1, 2, \dots$

③ If $\|v^{n+1} - v^n\| < \frac{\epsilon}{2\lambda} (1-\lambda)$ then stop and return
 v^{n+1} and d^* , otherwise continue to step 2

Let d attain max, then

$$P_d = P_d^L + P_d^U \quad \begin{array}{l} \text{upper diagonal} \\ \text{including diagonals} \\ \text{lower triangular} \end{array}$$

Step 2 says: $v^{n+1} = (I - \lambda P_d^L)^{-1} r_d + (I - \lambda P_d^L)^{-1} \lambda P_d^U v^n$

Show: $v^{n+1} = (I - \lambda P_d^L)^{-1} r_d + (I - \lambda P_d^L)^{-1} \lambda P_d^U v^n$

Ans: as d attains max $v^{n+1} = r_d + \lambda P_d^L v^{n+1} + \lambda P_d^U v^n$, then follows

so, $I - \lambda P_d = Q_d - R_d$, and $Q_d^{-1} > 0$ and $R_d \geq 0$ then this is called regular splitting

In value iteration (V.I) $Q_d = I$ $R_d = \lambda P_d$

In Gauss Seidel V.I $Q_d = I - \lambda P_d^L$
 $R_d = \lambda P_d^U$

$$v^{n+1} = Q_d^{-1} r_d + Q_d^{-1} R_d v^n \quad (\text{G.S V.I})$$

Theorem: let (Q_d, R_d) be a regular splitting of $I - \lambda P_d$, $\forall d \in D^{\text{MD}}$, assume

$$\alpha = \sup_d \|Q_d^{-1} R_d\| < 1$$

a) For any $v^0 \in V$, the iterative scheme

$$\text{to } v_\lambda^* \quad v^{n+1} = \max_d \{ Q_d^{-1} r_d + Q_d^{-1} R_d v^n \} = T v^n, \text{ converges}$$

b) v_λ^* is unique fixed point of T

c) linear convergence, with rate $\leq \alpha$
 $\text{AARC} \leq \alpha$

converges $O(\beta^n)$ with $\beta \leq \alpha$

Proof: T is a contraction is already seen in problem sheet

$\|Tu - Tv\| \leq \alpha \|u - v\|$ (also seen in problem sheet)

$\Rightarrow \exists$ unique v^* s.t $Tv^* = v^*$ (Banach fixed point theorem)

sufficient to show $v^* = v_\lambda^*$

$$v^* = T v^*$$

$$\text{as } v^* \geq Q_d^{-1} r_d + Q_d^{-1} R_d v^*, \forall d \quad (\text{By defn of max})$$

$$\Rightarrow (I - Q_d^{-1} R_d) v^* \geq Q_d^{-1} r_d$$

$$\Rightarrow v^* \geq (I - Q_d^{-1} R_d)^{-1} Q_d^{-1} r_d \quad \left(\because \sup_d \|Q_d^{-1} R_d\| < 1 \text{ so } (I - Q_d^{-1} R_d)^{-1} \text{ is all positive and invertible} \right)$$

$$\Rightarrow v^* \geq \underbrace{(Q_d)(I - Q_d^{-1} R_d)^{-1}}_{Q_d - R_d} r_d$$

$$A^{-1} B^{-1} = (AB)^{-1}$$

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

$$\Rightarrow v^* \geq v_\lambda^\pi \text{ for } \pi = (d, d, \dots) \in \Pi^{\text{SD}}$$

$$\Rightarrow v^* \geq \sup_{\pi \in \Pi^{\text{SD}}} v_\lambda^\pi = v_\lambda^*$$

$$\Rightarrow v^* \geq v_\lambda^*, \text{ now as } \exists d^* \text{ s.t } v^* = r_{d^*} + \lambda P_{d^*} v^* \quad v^* = L_{d^*} v^* \Rightarrow v^* = v_\lambda^{(d^*, d^*)}$$

as Ld^* has a unique fixed point
now, $\exists d^* \text{ s.t } v^* = v_{\lambda}^{(d^*, d^*, \dots)}$

$$\Rightarrow v^* = v_{\lambda}^*$$

part (c) follows from previous theorem

Note: (Q_1, R_1) and (Q_2, R_2) two regular splittings of $I - \lambda P$ if $R_2 \leq R_1 \leq \lambda P$
then $\|Q_2^{-1}R_2\| \leq \|Q_1^{-1}R_1\|$

$$\text{then } R_1 = \lambda P_d$$

$$R_2 = \lambda P_d u$$

then $R_2 \leq R_1$ holds as we are just removing some values,

$$\Rightarrow \|(I - \lambda P_d^L)^{-1} \lambda P_d v\| \leq \|\lambda P_d\| = \lambda < 1$$

$$\Rightarrow \sup_d \|(I - \lambda P_d^L)^{-1} \lambda P_d v\| \leq \|\lambda P_d\| = \lambda < 1$$

$$\Rightarrow \alpha \leq \lambda < 1$$

Note: If $P_d \neq P_d^v$ then $\alpha < \lambda$

so, to recap, we saw $v_{\lambda}^{\pi}(s) = \mathbb{E}_s^{\pi} \left[\sum_{t=1}^{\infty} \lambda^{t-1} r(x_t, y_t) \right]$ our objective function

① suffices to look at markovian randomised policies

② $L_d v = r_d + \lambda P_d v$, we saw L_d has unique fixed point
 $v_{\lambda}^{(d, d, \dots)}$

③ $L v = v$, if $v \geq L v$ then $v \geq v_{\lambda}^*$
 $v \leq L v$ then $v \leq v_{\lambda}^*$
 $v = L v$ then $v = v_{\lambda}^*$

④ L is a contraction, by banach fixed point, \exists unique v^*
s.t $L v^* = v^*$

⑤ Converging: $d^* \in \arg \max \{r_d + \lambda P_d v^* \}$

⑥ V.I

Exe: Show the \leq part of (b)

$$\text{Ans: } L v_N = v_{N+1} \leq v_N$$

$$\Rightarrow L(v_{N+1}) \leq L(v_N)$$

:

$$\Rightarrow L^m(v_{N+m}) \leq L^m(v_N)$$

$$\Rightarrow v_{N+m+1} \leq v_{N+m} \forall m \geq 0$$

26th Sept:

$s, a, r(s, a), \lambda$ considering infinite horizon $x_1 \rightarrow y_1 \rightarrow x_2 \rightarrow y_2 \dots$ realization

$$\underbrace{x_t \text{ state} \rightarrow y_t \text{ action}}_{R_t = \text{random reward}}$$

$$R_t = r(x_t, y_t)$$

Note: $R_t = r(x_t, y_t, x_{t+1})$

$$r(s, a) = \sum_{s' \in S} P(s' | s, a) r(s, a, s')$$

\uparrow depends on future state

currently we look at: $\mathbb{E}_s^{\pi} [\sum_{t=1}^{\infty} \lambda R_t]$ called risk sensitive mdp

$d_s(a) \rightarrow \text{MD}$
 $\pi \leftarrow$ choose a particular action for every s d is for every s , \exists one action

$r_d = (r_{d(s)})_s = (r(s, d(s))) \forall s \rightarrow \text{vector}$
similarly P_d is a matrix

if $\pi = (d_1, d_2, \dots)$ we call it SMD or $d^\infty = \pi$

Note: we have already seen that its same to work with SMD

Policy evaluation:

$\pi = (d_1, d_2, \dots)$ $v^\pi = \text{vector s.t}$

$$v^\pi = (v_s^\pi, \forall s)$$

$$\text{if } \pi = d^\infty: v^\pi = r_{d_1} + \lambda P_d v^\pi \quad \tilde{\pi} = (d_2, d_3, \dots)$$

$$v^\pi = r_{d_1} + \lambda P_d v^\pi$$

also, $v^* = \max_d (r_d + \lambda P_d v^*)$ we assume max as everything else is finite

Note: To solve a very long N finite horizon MDP, we can convert that to $N \rightarrow \infty$ and solve, given λ is given

d is a decision rule, d^∞ is a policy (SMD), value of this:

$$v^d = r_{d_1} + \lambda P_d v^d$$

$$\Rightarrow v^d = (I - \lambda P_d)^{-1} r_d \rightarrow \text{value under } d^\infty$$

we want to improve v^d :

$$A^d \triangleq \arg \max_d (r_d + \lambda P_d v^d)$$

i.e. if $d_1 \in A^d$ then

\downarrow
 set of all decision rules $d \in A^d$ $d \notin A^d$

$$d_1(s) \in \arg \max_{a_s \in A_s} (r(s, a_s) + \lambda \sum_{s' \in S} P(s'|s, a_s) v^{d(s')})$$

case I: $d \in t^d$ then $v^d = v^*$ ^{optimal} of objective function

case II: $d \notin t^d$, then d is not optimal

Let's pick $d_i \in t^d$ s.t.

if $d(s) \in t^d(S)$ for some s

then $d_i(s) = d(s)$ else we pick

now we have d_i , we get v^{d_i} under d_i

so, $v^{d_i} > v^d \rightarrow$ improvement step

(at least one component

strictly better)

Now, as $t^{d_1} = \text{argmax}_d (r_d + \lambda P_d v^{d_1})$

$d_2 \in t^{d_1}$ if $d_1 \notin t^{d_1}$

\vdots
 $d_K \in t^{d_{K-1}}$ if $d_{K-1} \notin t^{d_K}$

algorithm will stop in finite steps as finite d (SMD) and it cannot repeat as improving, never get back discarded d_j

Lemma: $v^{d_1} > v^d$

Proof:

$$d_K \rightarrow v^{d_K}$$

$$d_{K+1} \in t^{d_K}$$

$v^{d_{K+1}} > v^{d_K} \rightarrow$ we want to show this

By defn of d_{K+1} :

$$r_{d_{K+1}} + \lambda P_{d_{K+1}} v^{d_K} > r_d + \lambda P_d v^{d_K} \quad \forall d$$

for $d = d_K$

$$\Rightarrow r_{d_{K+1}} + \lambda P_{d_{K+1}} v^{d_K} > r_{d_K} + \lambda P_{d_K} v^{d_K} = v^{d_K}$$

so we have to show, $r_{d_{K+1}} + \lambda P_{d_{K+1}} v^{d_K} \leq v^{d_{K+1}}$ and

$$\text{so, } v^{d_K} \leq v^{d_{K+1}}$$

now, $r_{d_{K+1}} + \lambda P_{d_{K+1}} v^{d_K} \geq v^{d_K}$

$$\Rightarrow r_{d_{K+1}} > (I - \lambda P_{d_{K+1}}) v^{d_K}$$

$$\Rightarrow (I - \lambda P_{d_{K+1}}) v^{d_{K+1}} > (I - \lambda P_{d_{K+1}}) v^{d_K}$$

$$\text{as } r_{d_{K+1}} = (I - \lambda P_{d_{K+1}}) v^{d_{K+1}}$$

$\leftarrow Ax \geq Ax' \text{ (true as, } \lambda > 0 \text{ & } A \text{ is positive)}$

now if A is positive matrix, $\Rightarrow Ax \geq Ax'$

$$\text{so, } v^{d_{K+1}} \geq v^{d_K}$$

Now, $v^d = r_d + \lambda P_d v^d$ is known to us, so let's consider any $v \in V$

s.t. $v \geq r_d + \lambda P_d v \quad \forall d$ component wise version

$$\text{i.e. } v(s) \geq r(s,a) + \lambda \sum_{s'} p(s'|s,a) v(s') \quad \forall s, a$$

new if $\nabla = \{v \mid v(s) \geq r(s,a) + \lambda \sum_{s'} p(s'|s,a)v(s'), \forall s,a\}$

then minimizer of $\nabla = v$ s.t v is optimal

$$\text{and } v(s) = r(s,a) + \lambda \sum_{s'} p(s'|s,a)v(s'), \forall s,a$$

this will always happen from $v_{k+1} \geq v_k$ argument (\because finite many)

Note: Our problem becomes: given r, p, λ :

for all $\alpha(s) \geq 0, \forall s, \sum \alpha(s) = 1$

$$\min_{v} \sum_s \alpha(s)v(s) \text{ is one solution s.t } \left\{ \begin{array}{l} v \text{ is: } v(s) \geq r(s,a) + \lambda \sum_{s'} p(s'|s,a)v(s') \quad \forall a,s \\ \text{primal LP} \end{array} \right.$$

This becomes LP as $(I - \lambda P_d)v \geq r_d$ is what we are solving, i.e

$$(1 - \lambda \sum_{s'} p(s'|s,a))v(s) \geq r(s,a) \quad \forall s,a$$

$$\max_{(\chi(s,a))} \sum_{s,a} \chi(s,a)r(s,a) \text{ where } \left\{ \begin{array}{l} \sum_a \chi(s,a) = 1 \\ \chi(s,a) \geq 0 \quad \forall s,a \\ \sum_a \chi(s',a) - \lambda \sum_{s'} p(s'|s,a)\chi(s,a) = r(s'), \quad \forall s' \\ \text{dual LP} \end{array} \right.$$

Now if v^* is solution for primal, then there will be a dual solution χ^*
(by duality theory $\sum \alpha(s)v^*(s) = \sum a \chi^*(s,a)$)

$$\xrightarrow{\chi^*} d^*(s,a) \triangleq \underline{\chi^*(s,a)} \quad \forall a, \forall s$$

$\sum_{a'} \chi^*(s,a')$ ↗ if we choose basis solution, then $d^* \in \mathbb{R}^A$

we claim $d^{*\infty}$ is optimal

$v^* \rightarrow$ value v of MDP

$(d^*)^\infty \rightarrow$ optimal policy for MDP

Note: $E_\alpha^\pi [\sum \lambda^{t-1} R_t]$

$x_1 \sim \alpha$ ↗ starting state probability of states

unconstrained MDP problem

now let, $J_{\alpha,\lambda}^\pi = E_\alpha^\pi [\sum \lambda^{t-1} R_t]$ where $v_{\alpha,\lambda} \triangleq \sup_\pi J_{\alpha,\lambda}^\pi$

let, $c_{\alpha,\lambda}^\pi \triangleq E_\alpha^\pi [\sum \lambda^{t-1} c_t]$, constraint optimisation problem:
↑ constraint function

$$c_t = c(x_t, y_t)$$

$$v_{\alpha,c} \triangleq \sup_\pi J_{\alpha,c}^\pi, \text{ s.t } c_{\alpha,\lambda}^\pi \leq \beta \text{ is a constraint}$$

this can be solved by using a constraint from dual

$$\sum c(s,a)\chi(s,a) \leq \beta$$

we get this from $d_{\underline{\chi}}(s,a) = \underline{\chi(s,a)} \rightarrow d_{\underline{\chi}}^\infty$

$$J_{\underline{\chi}}^\infty = \sum r(s,a)\chi(s,a) \text{ is a lemma}$$

$d_{\underline{\chi}}^\infty = \sum c(s,a)\chi(s,a)$ for $\underline{\chi}$ is dual is a policy for MDP