

Cs 747 foundations of intelligent and learning agents

G. Akash Reddy



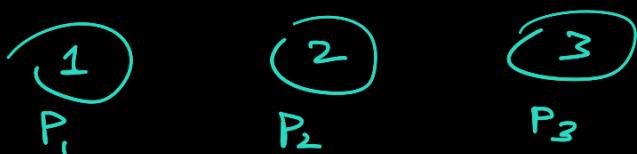
3 coding As - $3 \times 15 = 45$

midsem- 20

endsem- 35

\rightarrow Multi-armed bandits:-

coins example...



Q: get more heads, in 20 tosses.

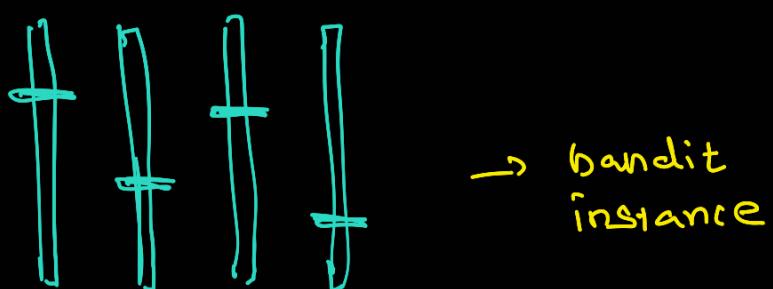
explore-exploit balance---

- explore coins to estimate P_1, P_2, P_3
- exploit the highest bias one.

Stochastic multi-armed bandit:-

monte-carlo tree search

What's this?



- n-arms each a Bernoulli RV.
- A be set of arms. ACA has mean p_A .
- let p^* be highest mean.

\rightarrow Algorithm:-

in case of bandits:

for $t = 0, 1, 2, 3, \dots T-1$:-

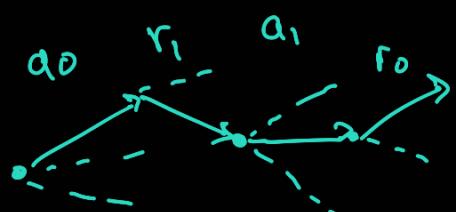
- say history $h^t = (a^0, r^0; a^1, r^1; a^2, r^2; \dots a^{t-1}, r^{t-1})$
- pick an arm a^t to "pull." $a^t \rightarrow$ arm no. for $t=1$
- obtain reward r^t . $r^t \in \{0, 1\}$

↳ reward for $t=1$

- * $T \rightarrow$ budget or "horizon" of experiment
- * Formally, a deterministic algorithm is a mapping
 - from set of all histories
 - to set of all arms.

(i.e. based on history, decides which arm to pull deterministically)
- * Formally, a randomised algorithm is a mapping
 - from set of all histories
 - to set of "all probability distributions over arms".

(i.e. based on history; assigns probability to pull, for each arm).
- * deterministic algorithm $\rightarrow 2^T$ histories possible.



$$h^T = (a^0, r^0, a_1, r_1, \dots, a_{t-1}, r_{t-1})$$

$$P(h^T) = \prod_{t=0}^{T-1} P(a^t | h^{t-1}) \cdot P(r^t | a^t)$$

~~~~~

algorithm      bandit instance

→  $\epsilon$ -greedy algorithms :-

$\epsilon \in [0,1]$  controls the amount of exploration.

$\epsilon G_1$  :

$t < \epsilon T$ ; sample arms randomly

at  $t = \lfloor \epsilon T \rfloor$  identify  $a^{\text{best}}$  with highest mean.

$t > \lfloor \epsilon T \rfloor$  pull  $a^{\text{best}}$

$\epsilon G_2$  :

if  $t < \epsilon T$ ; sample randomly

if  $t > \epsilon T$ ; pull highest empirical mean.

$\epsilon G_3$  :

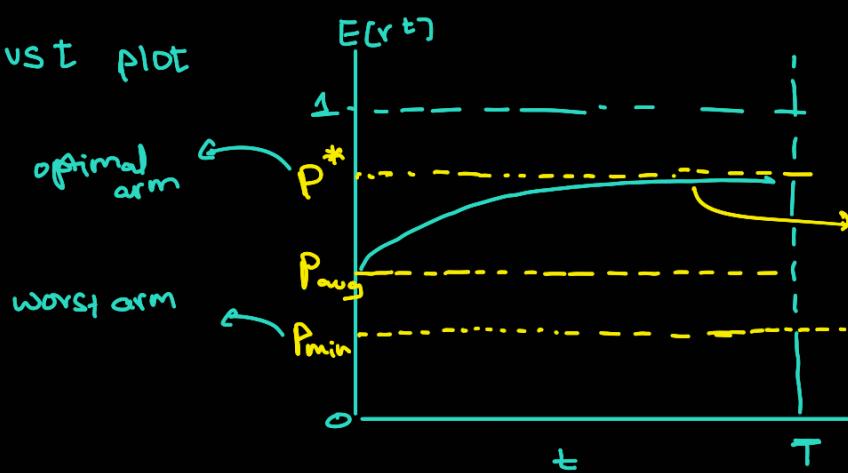
probability  $\epsilon \rightarrow$  pull randomly any arm

probability  $(1-\epsilon) \rightarrow$  pull highest mean one.

\* Does  $\epsilon G_1$  perform better than  $\epsilon G_2$ ? What is a "better" algorithm?

→ Visualizing performance:-

$E[r^t]$  vs  $t$  plot



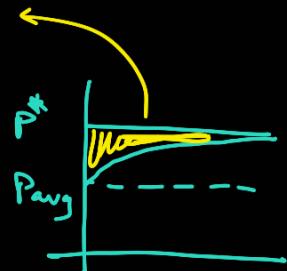
graph for a reasonable learning algorithm.

\* consider a arm; which pulls each arm at random.

Hence  $E[r^t] = \text{avg}\{P_1, P_2, \dots, P_n\}$  for any  $t \leq T$

- \* Maximum reward for period  $T$  is  $Tp^*$ .
- \* Actual expected reward =  $\sum_{t=0}^{T-1} E(r^t)$
- \* (Expected cumulative) Regret of algorithm for horizon  $T$  is

$$R_T = Tp^* - \sum_{t=0}^{T-1} E(r^t)$$



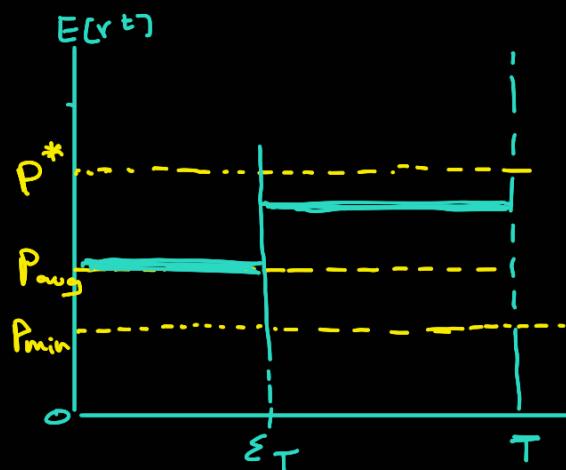
- \* We want  $R_T$  to be smaller. In fact  $T \rightarrow \infty \frac{R_T}{T} = 0$ .

"Sub-linear" regret

We want algorithms which give sub-linear regret.

"don't all of them do that?"

→ EG<sub>1</sub> algo:



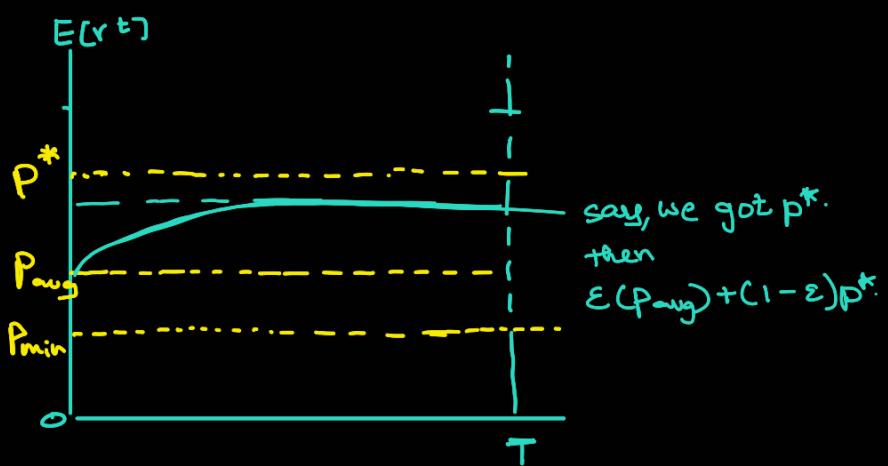
$$\therefore R_T = (P^* - P_{\text{avg}})\epsilon_T + \text{some value}$$

$$\therefore \lim_{T \rightarrow \infty} R_T = (P^* - P_{\text{avg}})\epsilon \neq 0$$

Not sub-linear. | Same analysis for EG<sub>2</sub>

→ EG3:

uniformly with  $\varepsilon$   
exploit with  $(1-\varepsilon)$



$$\therefore \frac{R_T}{T} \geq p^* - [\varepsilon(p_{avg}) + (1-\varepsilon)(p^*)]$$

$$\geq \varepsilon(p^* - p_{avg})$$

$$\therefore \lim_{T \rightarrow \infty} \frac{R_T}{T} \geq \underline{\varepsilon(p^* - p_{avg})}$$

linear regret.  
Not sub-linear regret.

→ How to achieve sub-linear regret?

\* Two conditions - C<sub>1</sub> and C<sub>2</sub>

\* C<sub>1</sub>: Infinite Exploration

\* In the limit  $T \rightarrow \infty$ , each arm must be pulled an  $\infty$  no. of times.

Explanation: ➤ Say an arm is at max pulled ' $U$ ' times,  
➤ with a (slight) probability of  $(1-p^*)^U$ ; the most optimal arm will have empirical mean 0.

∴ The Regret  $\frac{R_T}{T} \geq \underbrace{(1-p^*)^U}_{\text{not going to 0 with } T \rightarrow \infty} (p^* - p_2)$

\* to avoid this slightly probable misfortune; we need to explore endlessly.

## \* C<sub>2</sub>: Greed in the limit

\* Also have to exploit enough. let exploit(T) denote no. of pulls that are greedy wrt. empirical mean...

$$\lim_{T \rightarrow \infty} \frac{E[\text{exploit}(T)]}{T} = 1$$

\* let  $\bar{\mathcal{I}}$  be the set of all bandit instances with reward-means strictly less than 1,

### RESULT:

An algorithm L achieves sub-linear regret on all instances  $I \in \bar{\mathcal{I}}$  if and only if it has C<sub>1</sub> and C<sub>2</sub> properties on all  $I \in \bar{\mathcal{I}}$ .

GILIE conditions, greedy in limit & infinite exploration.

$Eg_1, Eg_3 \rightarrow$  both have infinite exploration.

$\rightarrow$  don't have greedy in the limit.

(both are  $(1-\varepsilon)$  in limit)

$\rightarrow$  Making our  $Eg_1$  as sub-linear regret:-

\* say  $\epsilon_T = \frac{1}{\sqrt{T}}$ . nice! ( $\epsilon_T$ -first algorithm)

then  $\epsilon_T \cdot T = \sqrt{T} \rightarrow$  infinite times exploration

$$\& \frac{\text{Exploit}(T)}{T} = 1 - \frac{1}{\sqrt{T}} \xrightarrow{\text{GIL.}} 1 \text{ as } T \rightarrow \infty$$

\* For  $\epsilon_t$ -greedy

here, say  $\epsilon_t = \frac{1}{1+t}$  :  $\frac{1}{(1+t)^2}$  won't work! IE fails I think.

∴ at time  $t$ :  $\epsilon_t$  exploration  $\approx 1-\epsilon_t$  exploitation.

\* Checking C<sub>1</sub>: each arm gets  $\sum_{t=0}^{T-1} \frac{1}{n(1+t)} = \Theta\left(\frac{\log T}{n}\right)$  pulls... infinite!

\* Checking C<sub>2</sub>: since exploration is  $\approx \log(T)$

$$\frac{\text{exploit}(t)}{T} = 1 - \frac{\log(t)}{T} \rightarrow 1 \text{ as } T \rightarrow \infty$$

∴ GIL holds!

→ Lower bound on regret :-

Lai and Robbins lower bound on regret

Let L be an algorithm such that for every  $I \in \bar{\mathcal{I}}$  and for every  $\alpha > 0$ , as  $T \rightarrow \infty$

$$R_T(L, I) = O(T^\alpha) \text{ i.e. sub-polynomial}$$

$\therefore \alpha > 0$

Then for every instance  $I \in \bar{\mathcal{I}}$ , as  $T \rightarrow \infty$ :

$$\frac{R_T(L, I)}{\ln(T)} \geq \sum_{a: P_a \neq P^*} \frac{P^*(I) - P_a(I)}{KL(P_a, P^*)}$$

where for  $x, y \in [0, 1]$

$$KL(x, y) \stackrel{\text{def}}{=} x \ln \frac{x}{y} + (1-x) \ln \frac{1-x}{1-y}$$

→ UCB - upper confidence bounds sampling :-

- \* at time  $t$ ; define  $UCB_a^t = \hat{P}_a^t + \sqrt{\frac{2\ln t}{u_a^t}}$
- pull each arm once first.
- \* pull the arm with the highest UCB value.
- \* Achieves regret  $O(\log(T))$ . optimal dependency.  
but factor is not exactly optimal.

→ KLUCB:

- \* Identical to UCB algorithm; but definition is different!

$$UCB-KL_a^t = \max_{q \in [\hat{P}_a^t, 1]} \text{st. } u_a^t \cdot KL(\hat{P}_a^t, q) \leq \ln t + c \cdot \ln(\ln t)$$

$c \geq 3$

- pull arm with max KL-ucb value.

$$KL(x, y) \stackrel{\text{def}}{=} x \ln \frac{x}{y} + (1-x) \ln \frac{1-x}{1-y}$$

- \* KL-ucb is tighter bound than UCB.

KL-ucb asymptotically matches loi-robbins bound.

nice!

so, bandits is completely solved...  
at least this particular model.

## → Thompson Sampling :-

- \* Beta distributions :-

$$\boxed{\beta(x; \alpha, \beta) = \frac{(\alpha+\beta-1)!}{(\alpha-1)! (\beta-1)!} x^{\alpha-1} (1-x)^{\beta-1}} \quad x \in [0, 1]$$

$$\text{Variance} = \frac{\alpha\beta}{(\alpha+\beta)^2 (\alpha+\beta+1)} \quad \text{Mean} = \frac{\alpha}{\alpha+\beta}$$

- \*  $\beta(1, 1)$  is uniform distribution.

- \* At time  $t$ , let arm  $a$  have  $s_a^t$  successes &  $f_a^t$  failures.

$\beta(s_a^t + 1, f_a^t + 1)$  represents the belief or mean of arm  $a$  at time  $t$ .

- \* Sample a value for arm  $a$  from this distribution.

- \* Pull the arm with the highest sampling value.

- \* Achieves optimal regret; excellent in practice.

Lai & Robbins lower bound.

## $\rightarrow$ Hoeffding's inequality

$X \rightarrow$  random variable  $\{0, 1\}$ ;  $E(X) = \mu$

$x_1, x_2 \dots x_n$  are iid of  $X$ .

\*  $\bar{x} = \frac{1}{n} \sum x_i$

then

$$P(\bar{x} > \mu + \varepsilon) \leq e^{-n\mu\varepsilon^2}$$

$$P(\bar{x} < \mu - \varepsilon) \leq e^{-n\mu\varepsilon^2}$$

Applications  
in slides

## $\rightarrow$ KL-inequality :-

\* Same setting as before.

$$P\{\bar{x} > \mu + \varepsilon\} \leq e^{-n\mu K_L(\mu + \varepsilon, \mu)}$$

$$P\{\bar{x} < \mu - \varepsilon\} \leq e^{-n\mu K_L(\mu - \varepsilon, \mu)}$$

tighter than  
Hoeffding's  
inequality.

\* Using above results; proof that UCB is  $O(1/\log T)$ .  
*discussed in slides*

*asked in  
midsem.. { "Analysis of UCB"*

\* Analysis of Thompson Sampling has basics in Bayesian inference.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

\* the logic to select which arm to pull is justified as

"The Beta distributions reflect our beliefs about the arms. Collectively, they represent our beliefs about the bandit instance.

We are sampling a bandit instance & pulling optimal arm wrt. that instance".

(nice one!)

## → Markov Decision Problems MDPs :-

\* Elements of an MDP  $M = (S, A, T, R, \gamma)$

$S$  : states

$A$  : Actions

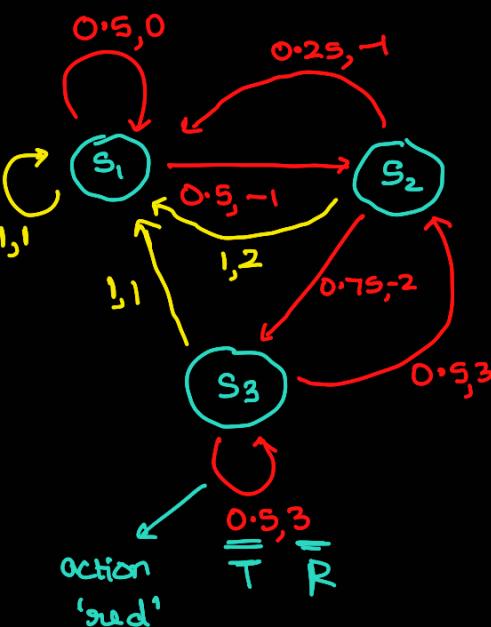
$T$  : transition probabilities

$R$  : rewards

$\gamma$  : discount factor--?  
 $\in [0, 1]$

$$T(s, a, s') \in [0, 1]$$

$$R(s, a, s') \in [-R_{\max}, R_{\max}]$$



$$\text{Here } S = \{S_1, S_2, S_3\}$$

$$A = \{\text{'red'}, \text{'yellow'}\}$$

$T, S$  are on every edge.

\* How does agent pick  $a^t$ ?

In principle, it can decide by looking at the history:

$$s^0, a^0, r^0, s^1, a^1, r^1, s^2, \dots$$

For now, assume  $a^t$  depends only on  $s^t$ . (markovian!)

stationary  
deterministic

(↳ not stochastic)

single action  
leads to  
single state;  
not probability  
over states.

→ State values for policy  $\pi$ :

$$\text{for } s \in S, V^\pi(s) \stackrel{\text{def}}{=} E_\pi [r^0 + r^1 + r^2 + \dots], S^0 = S$$

But, this sum goes  $\infty$   
 might

Hence, we discounting factor

$$V^\pi(s) \stackrel{\text{def}}{=} E_\pi [r^0 + \gamma r^1 + \gamma^2 r^2 + \gamma^3 r^3 + \dots] \quad S^0 = S$$

Larger the  $\gamma$ , farther the "lookahead" while considering reward.

→ MDP planning:-

Given  $M = (S, A, T, R, \gamma)$ ; find a policy  $\pi^*$  from the set of all policies  $\Pi$  such that  $\forall s \in S; \forall \pi \in \Pi$

$$V^{\pi}(s) \geq V^{\pi^*}(s)$$

- \* Every MDP is guaranteed to have atleast one such markovian, deterministic, stationary optimal policy.
- \* Can have more than one optimal policy but only one unique optimal value function.

→ Policy evaluation: Finding  $V^\pi(s)$  from  $\pi$ :-

$$V^\pi(s) = E_\pi (r + \gamma r_1 + \gamma^2 r_2 + \dots)$$

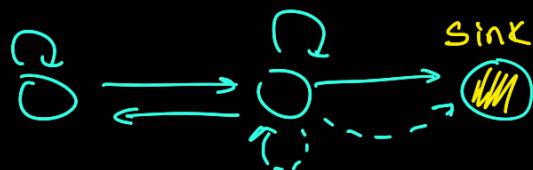
$$V^\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') (R(s, \pi(s), s') + \gamma V^\pi(s'))$$

Bellman equations

now we have 'n' such equations  
& 'n' unknowns

- \* Guaranteed to have unique solution if  $\gamma < 1$ .
- \* Using this, we iterate over all policies to find the optimal policy.  
Exponential!
- lets see better algorithms.

→ Episodic tasks :-



we considered infinitely long tasks.

- \* Now; episodic tasks have special sink states from which there are no outgoing tasks or rewards.

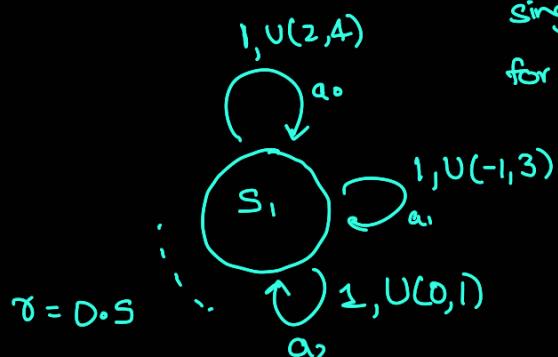
\* Additionally, from every non-terminal state, for every policy, there is non-zero probability of reaching the sink in finite time.

\* here; we can talk about net rewards; without discounting future rewards

$$V^\pi(s) \triangleq \mathbb{E}_{\pi}[r_0 + r_1 + \dots] \mid s^0 = s$$

WILL be finite for episodic tasks.

↳ A familiar MDP?



single state;  $k$ -actions

for  $a \in A$ , treat  $R(s, a, s')$  as a random variable

such an mdp → multi armed bandit.

## → Complete, normed vector Spaces :-

- \* A norm  $\| \cdot \|$  associates a length with each vector.
- \* A complete, normed vector space -  $(X, \| \cdot \|)$  is one in which every Cauchy sequence has a limit in  $X$ .  
no need to contemplate
- \* A complete, normed vector space is called a Banach space.

→ Definitions: say  $(X, \|\cdot\|)$  is a normed vector space s.t.  $0 \leq l < 1$ .

contraction mapping: A mapping  $Z: X \rightarrow X$  is called a contraction mapping with contraction factor  $l$ ; if  $\forall u, v \in X$

$$\|Zu - Zv\| \leq l \|v - u\|$$

Fixed point:  $x^* \in X$  is called fixed point of  $Z$  if

$$Zx^* = x^*$$

→ Banach's fixed point theorem:-

- \*  $(X, \|\cdot\|)$  be a banach space
- \*  $Z: X \rightarrow X$  be contraction mapping with  $l \in [0, 1)$

1.  $Z$  has unique fixed point  $x^* \in X$
2. for  $x \in X$ ;  $m \geq 0$   $\|Z^m x - x^*\| \leq l^m \|x - x^*\|$

"finally; all points collapse onto fixed point".

→ Bellmann Optimality operator-

- \*  $S = \{s_1, s_2, s_3, \dots\}$  states of an MDP.

- \*  $F: S \rightarrow \mathbb{R}$  could be MDP value function  
is a point in  $\mathbb{R}^n$ . (fine!)

$B^*: \mathbb{R}^n \rightarrow \mathbb{R}^n$  for MDP  $(S, A, T, R, \gamma)$  is

Bellmann optimality operator

point in  $\mathbb{R}^n$ ;  $S \rightarrow \mathbb{R}$ .

$$(B^*(F))(s) \stackrel{\text{def}}{=} \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{R(s, a, s') + \gamma F(s')\}$$

- \*  $\|F\|_\infty = \max \{ |f_1|, |f_2|, \dots, |f_n|\}$  this is called max norm  $\|.\|_\infty$
- \*  $(\mathbb{R}^n, \|.\|_\infty)$  is a Banach space. Established!

Fact:

$B^*$  is contraction mapping in  $(\mathbb{R}^n, \|.\|_\infty)$  Banach space with contraction factor  $\gamma$ .

\* Proving  $B^*$  is a contraction mapping:-

we use

$$\left| \max_a f(a) - \max_a g(a) \right| \leq \max_a |f(a) - g(a)|$$

\* we can show  $\|B^*(F)(s) - B^*(G)(s)\|_\infty \leq \gamma \|F(s) - G(s)\|_\infty$

→ The fixed point of  $B^*$ :

\* there is a unique fixed point of  $B^*$ .

Denote the fixed point as  $v^*: S \rightarrow \mathbb{R}$  & so  $B^*(v^*) = v^*$

→ for  $s \in S$

$$v^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') (R(s, a, s') + \gamma v^*(s'))$$

//

$n$  variables,  $n$  equations but  
non-linear

these are  
Bellmann  
optimality  
equations.

- |                       |   |                                            |
|-----------------------|---|--------------------------------------------|
| 1. Value Iteration    | } | 3 family of algorithms<br>to compute $v^*$ |
| 2. Linear Programming |   |                                            |
| 3. Policy Iteration   |   |                                            |

## Fact:

$V^*(s)$  is the value function of every policy  $\pi^*: S \rightarrow A$  such that for  $s \in S$

$$\pi^*(s) = \arg\max_a \sum_{s' \in S} T(s, a, s') (R(s, a, s') + \gamma V^*(s'))$$

## 1) Value iteration:-

- \* Iterative approach to compute  $V^*$ .

$$v_0 \xrightarrow{B^*} v_1 \xrightarrow{B^*} \dots v^* \quad (\text{Awesome lol!})$$

- \* Popular, easy to implement, quick to convergence

## 2- Linear Programming :-

- \* To solve real valued  $x_1, x_2, \dots, x_n$ ;
  - given linear function is maximised
  - given linear constraints are satisfied.
- \* Modern linear programming solvers can solve LPs with thousands / millions of variables / constraints in reasonable time (hours / days)
  - common approaches: simplex, interior-point methods
- \* Engineer's focus is on formulating, rather than solving LPs.

## MDP planning using linear programming:-

1) Bellmann optimality equations :  $s \in S$

$$v^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma v^*(s') \}$$

non linear.

we make it linear as:  
constraints

$$v(s) \geq \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma v^*(s') \}$$

$\forall s \in S$   
 $\forall a \in A$

$n \cdot k$  linear constraints!

\* Note that  $v^*$  is in the feasible set.

\* Now; construct an objective function; so that  $v^*$  is the sole optimizer.

Digression :-

Vector comparison :-

\* for  $x: S \rightarrow \mathbb{R}$  and  $y: S \rightarrow \mathbb{R}$ ; we define

$$x \geq y \Leftrightarrow \forall s \in S \quad x(s) \geq y(s)$$

$$x > y \Leftrightarrow x \geq y \text{ and } \exists s \in S: x(s) > y(s)$$

\* And naturally, for policies  $\pi_1, \pi_2$  we define

$$\pi_1 \geq \pi_2 \Leftrightarrow v^{\pi_1} \geq v^{\pi_2}$$

$$\pi_1 > \pi_2 \Leftrightarrow v^{\pi_1} > v^{\pi_2}$$

\* Note: we can have two incomparable policies.

$$\pi_1 \succ \pi_2 \text{ and } \pi_2 \succsim \pi_1 \iff \sqrt{\pi_1} = \sqrt{\pi_2}$$

\*  $B^*(\cdot)$  preserves  $\succeq$ :

Fact: for  $x: S \rightarrow \mathbb{R}$

$y: S \rightarrow \mathbb{R}$

$$x \succeq y \Rightarrow B^*(x) \succeq B^*(y)$$

Proof can be done -

$$\max_a f(a) - \max_a g(a) \geq \min_a (f(a) - g(a))$$

✓

\* Lets see the feasible set of our LP:-

Each  $v: S \rightarrow \mathbb{R}$  in our set satisfies  $v \succeq B^*(v)$   
 $\Rightarrow v \succeq B^{*2}(v) \succeq B^{*3}(v) \dots \succeq v^*$

∴ each  $v$  in the feasible set  $v \succeq v^*$

We linearise this as:

$$\sum_{s \in S} v(s) \geq \sum_{s \in S} v^*(s)$$

Hence; we need to minimize  $\sum_{s \in S} v(s)$ .

$$\Rightarrow \boxed{\text{maximize: } -\sum_{s \in S} v(s)}$$

subject to

our LP formulation is completed.

$$\boxed{v(s) \geq \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma v(s') \} \quad \forall s \in S \\ \forall a \in A}$$

→ Policy Iteration:-

→ Action value function:-

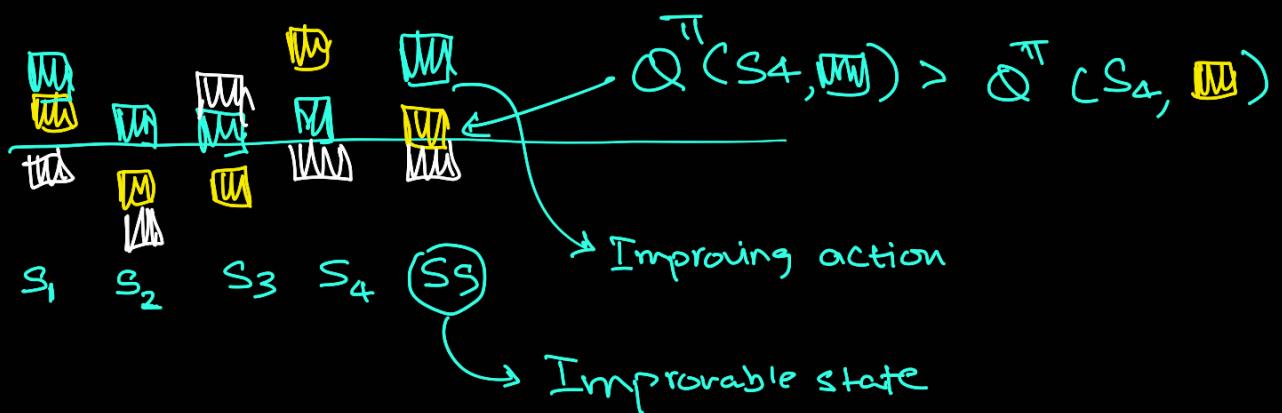
$$Q^\pi(s, a) \stackrel{\text{def}}{=} E[\gamma^0 + \gamma r^1 + \gamma^2 r^2 + \dots] \mid s^0 = s, a^0 = a, a^t = \pi(s^t) \text{ for } t \geq 1$$

$Q^\pi: S \times A \rightarrow \mathbb{R}$  is called action value function of  $\pi$ .

$$Q^\pi(s, a) = \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma V^\pi(s') \}$$

- \* All optimal policies have same optimal action value functions
  - expected reward, starting from state  $s$  & playing action  $a$ , & following policy  $\pi$  from next turn.

→ Policy Improvement:-



Policy Improvement:

Given  $\pi^0$

- pick one or more improvable states  $s_i$  in these
- Switch to an arbitrary improving action

$$\pi_{\text{new}} > \pi_{\text{old}}$$

## → Policy Improvement Theorem:-

$$V^\pi(s) \equiv Q^\pi(s, \pi(s))$$

\* for  $\pi \in \Pi, s \in S$ ;

$$IA(\pi, s) \stackrel{\text{def}}{=} \{a \in A : Q^\pi(s, a) > V^\pi(s)\}$$

\* for  $\pi \in \Pi$ ,

$$IS(\pi) \stackrel{\text{def}}{=} \{s \in S : |IA(\pi, s)| \geq 1\}$$

\* Suppose  $IS(\pi) \neq \emptyset$  &  $\pi' \in \Pi$  is obtained by policy improvement on  $\pi$ ;

### Policy Improvement Theorem:

- 1) if  $IS(\pi) = \emptyset$ ; then  $\pi$  is optimal. Else
- 2) if  $\pi'$  is obtained by policy improvement on  $\pi$ ,  
 $\pi' > \pi$

## → Bellman operator (No optimality):-

\* for policy  $\pi$ ; we define  $B^\pi: \mathbb{R}^n \rightarrow \mathbb{R}^n$  as:

$$(B^\pi(x))(s) \stackrel{\text{def}}{=} \sum_{\substack{s' \in S \\ \text{no max } a}} T(s, \underline{\pi}(s), s') \{ R(s, \underline{\pi}(s), s') + \gamma x(s') \}$$

### Facts:

1.  $B^\pi$  is a contraction mapping with  $\gamma$ . Like  $B^*$ .

2.  $\lim_{k \rightarrow \infty} (B^\pi)^k(x) = v^\pi$

3.  $x \geq y \Rightarrow B^\pi(x) \geq B^\pi(y)$

Note: for  $\pi, \pi' \in \Pi$  &  $s \in S$

$$B^{\pi'}(v^\pi)(s) = Q^\pi(s, \pi'(s))$$

### → Policy iteration Algorithm:-

```

 $\pi \leftarrow \text{arbitrary policy}$ 
while  $\pi$  has I.S.
   $\begin{cases} \pi' \leftarrow \text{policy improvement}(\pi) \\ \pi \leftarrow \pi' \end{cases}$ 
return  $\pi$ 

```

\* No. of iterations  
depend on the  
switching strategy.  
→ choose all I.S.

## → A more general class of policies:-

\* till now, we only saw

markovian, deterministic, stationary policies.

Why  
not

history  
dependent?

Stochastic?

non-stationary?

**FACT:** NO such non-markovian, stochastic, non-stationary policy  
can out-perform our  $V^*$

## → Different Policy iterations:-

1) Howard's policy iteration:-

Greedy → improve all improvable states

2) Random policy iteration:-

include/exclude a state → uniformly decided

3) Simple iteration policy:-

- Assume an indexing of states
- Improve highest index state

Upper & lower bounds -

- > Analysis
- > Basic tools
- > Open problems

} refer  
lecture  
slides

\* Now; what if  $T, R$  were not given, but had to be learned from interaction?  
Can we still learn to act optimally?

yes! the reinforcement learning problem.



→ Monte-Carlo methods:-

??

monte-carlo integration?

Reinforcement learning:-

1<sup>st</sup> class - intro : The before midsem ---  
2<sup>nd</sup> class - missed - Friday after midsem