

REINFORCEMENT

LEARNING :

MCTS + UCT

- RL is a trial and error learning but is based on long-term rewards.
- was used in Alpha zero v/s Stockfish chess engine

NOTE - MARK PROBLEM :

• Convex Evaluation:

1) Mid sem : 25

2) End sem : 25

3) 2 quizzes : 30%.

4) Mini-project : 20%.

$$P\{e_1, e_2\} = P\{e_1\} \underbrace{P\{e_2/e_1\}}_{(1)}$$

$$P\{e_2, e_1\} = P\{e_2\} P\{e_1/e_2\}$$

$$= Y_3$$

$$P\{T, e_1\} = P\{T\} P\{e_1/T\}$$

$$= (Y_3) \quad (Y_2)$$

$$= Y_6$$

$$P\{T, e_2\} = \overbrace{P\{T\}} P\{e_2/T\}$$

$$= Y_6$$

R.V:

RV is a fun. from

outcome space to \mathbb{R} .

(has a dist. associated with it)

- $E[x]$ is like centre of mass,
- $\text{var}[x]$ is like moment of inertia.

• If $v \sim \text{uniform}[0, 1]$ then,

$$w = -\ln(1-v) \sim \text{exp}$$

3 (7/1):

LAW OF LARGE NUMBERS:

• let x_1, x_2, \dots, x_n be IID R.V's

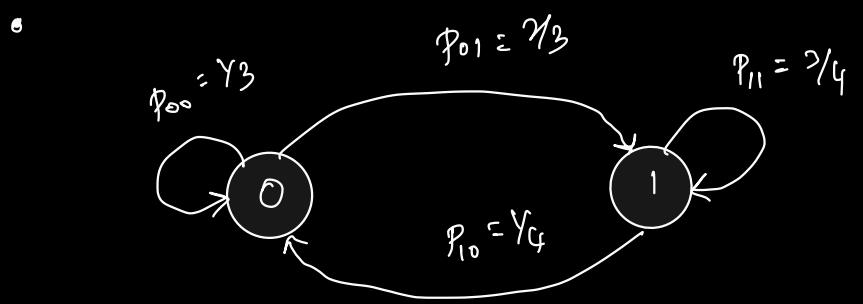
Then $s_n := \underbrace{\frac{1}{n} \sum_{i=1}^n x_i}_{\text{seq. of RV's}} \rightarrow E[x]$

MARKOV CHAIN:

• Consider, $\{x_n\}_{n \geq 0}$ seq. of RV's,

def: $\Pr \{ x_n = x \mid x_{n-1}, \dots, x_1 \} = \Pr \{ x_n = x \mid x_{n-1} \}$

• IID sequence is a trivial Markov chain



24:

Markov chain :

Characteristics:
Initial distribution : $\pi_0^{(n)} = P\{X_0 = x\}$

, Transition probabilities :

$$P\{X_{n+1} = j \mid X_n = i\} = P_{ij}$$

• Computing the probability of being in state 'i' at time 'n' :

wkT , $\begin{bmatrix} P^n \\ \downarrow \end{bmatrix}_{ij} = P\{X_n = j \mid X_0 = i\}$
PTM

$$\Rightarrow \pi_n^T = \pi_0^T P^n$$

• let N denote # of tosses before two consecutive heads.

$$N \in \{2, 3, \dots, \infty\}$$

$$E[N] = \sum_{n=1}^{\infty} n P\{N = n\}$$

$$P\{N=2\} = y_4 \quad \{HH\}$$

$$P\{N=3\} = y_8 \quad \{THH\}$$

$$P\{N=4\} = y_8 \quad \{HTHH, TTHH\}$$

Tedious Job

Alternative way:

$$\begin{aligned} \cdot E[N] &= E[N \mid \text{first toss} = T] P\{\text{First toss} = T\} \\ &\quad + E[N \mid \text{First toss} = H] P\{\text{First toss} = H\} \end{aligned}$$

$$\begin{aligned} \cdot E[N \mid \text{First toss} = T] &= \sum_{n=2}^{\infty} n P\{N \mid \text{1st toss} = T\} \\ &= \sum_{n=2}^{\infty} n P\{N \mid \text{1st toss} = T\} \\ &= \sum_{k=2}^{n=3} (k+1) P\{N \mid \text{1st toss} = T\} \\ &= \sum_{k=2}^{\infty} k P\{N \mid \text{1st toss} = T\} + \sum_{k=2}^{\infty} P\{N \mid \text{1st toss} = T\} \\ &= \underbrace{\sum_{k=2}^{\infty} k P\{N \mid \text{1st toss} = T\}}_{= E[N \mid \text{1st toss} = T]} + 1 \\ &= \frac{E[N] + 1}{-} \end{aligned}$$

$$\begin{aligned}
\mathbb{E}[N \mid 1^{\text{st}} \text{ toss} = H] &= \mathbb{E}[N \mid 1^{\text{st}} \text{ toss} = H, 2^{\text{nd}} \text{ toss} = H] \\
&\quad + \mathbb{E}[N \mid 1^{\text{st}} \text{ toss} = H, 2^{\text{nd}} \text{ toss} = T] \\
&\quad + \mathbb{E}[N \mid 1^{\text{st}} \text{ toss} = T, 2^{\text{nd}} \text{ toss} = H] \\
&= 2(\gamma_2) + (\mathbb{E}[N] + 2)\left(\frac{1}{2}\right) \\
&= 1 + (\mathbb{E}[N] + 2)\gamma_2 \\
&= \frac{\mathbb{E}[N]}{2} + 2
\end{aligned}$$

$$\Rightarrow \mathbb{E}[N] = (\mathbb{E}[N] + 1)\gamma_2 + \left(\frac{\mathbb{E}[N]}{2} + 2\right)\gamma_2$$

$$\mathbb{E}[N] = \mathbb{E}(N) \left(\frac{3}{4}\right) + \frac{3}{2}$$

$$\mathbb{E}[N] \left(\frac{1}{4}\right) = 3/2$$

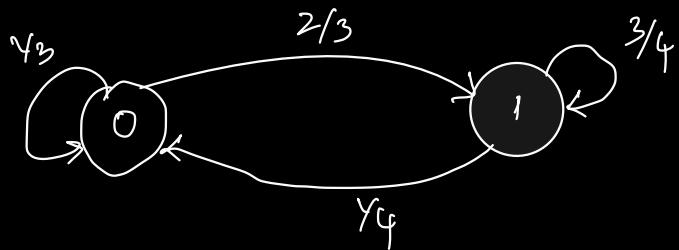
$$\mathbb{E}[N] = 6$$

$$\overbrace{=}^{<}$$

STATIONARY DISTRIBUTION :

$$\pi_n^T P = \pi_{n+1}^T$$

$$\pi^T P = \pi^T \quad \left. \begin{array}{l} \\ \end{array} \right\} \begin{array}{l} \text{Eigen. vector with} \\ \text{eigen value } = 1 \end{array}$$



$$\begin{bmatrix} x & y \end{bmatrix} \begin{Bmatrix} y_3 & 2/3 \\ y_4 & 3/4 \end{Bmatrix} = \begin{bmatrix} x & y \end{bmatrix}$$

$$\frac{x}{3} + \frac{y}{4} = x$$

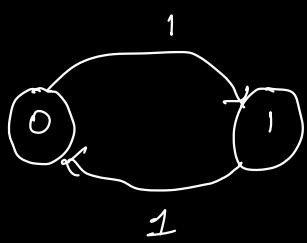
$$\frac{2}{3}x + \frac{3}{4}y = y \quad -8x + 3y = 0$$

$$x + y = 1$$

$$4x + 3y = 12x \quad -7x + 4y = 1$$

$$8x + 9y = 12x$$

Ex :



BANDITS : (stateless RL)

12/1:

, we look at Bernoulli Bandits with 2 arms and then we generalize it.

. For two arms, action space $A = \{1, 2\}$

. Let μ_i denote the avg. reward of arm i .

. Let reward at step ' t ' be R_t and A_t be the action at time ' t '.

. Our goal is to max $\sum_{t=1}^T R_t$.

Bnt here we don't have the expectation, so

it's like a realization.

. We choose A_t 's such that avg. reward is maximized i.e. max $\mathbb{E} \left[\sum_{t=1}^T R_t \right]$

. Notion of regret :

Regret = Max. avg. reward that can be obtained till T
- Total expected reward obtained by our algorithm.

Let optimal arm be i^* such that $\mu^* = \max_i \mu_1, \mu_2$

$$\Rightarrow \text{Regret} = \mu^* T - \underbrace{\mathbb{E} \left[\sum_{i=1}^T R_i \right]}_{\downarrow \text{Randomness w.r.t reward dist. and also the arm selection}}$$

Let N_1 : # of times arm 1 is pulled
 N_2 : # of times arm 2 is pulled

$$\Rightarrow \mathbb{E} \left[\sum_{i=1}^T R_i \right] = \mu_1 N_1 + \mu_2 N_2$$

$$\Rightarrow R_T = \frac{\mu^* T}{N_1 + N_2} - (\mu_1 N_1 + \mu_2 N_2)$$

$$R_T = (\mu^* - \mu_1) N_1 + (\mu^* - \mu_2) N_2$$

$$R_T = \Delta_1 N_1 + \Delta_2 N_2$$

Intuitive understanding
 $\Rightarrow \text{Regret} \propto \text{Avg. # of pulls of a sub-optimal arm}$

GREEDY POLICY:

For $t=1$ pull arm 1

For $t=2$ pull arm 2

For $t \geq 2$ pull the arm with highest estimated

avg - reward.

The estimates, $\hat{\mu}_1 = \frac{R_1^{(1)} + R_2^{(2)} + \dots + R^{(n_1)}}{n_1}$

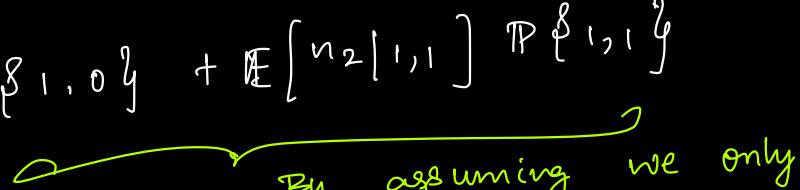
$$\hat{\mu}_2 = \frac{R_1^{(2)} + R_2^{(2)} + \dots + R^{(n_2)}}{n_2}$$

Regret $\geq O(T)$

To see why regret is lower-bounded by $O(T)$
 assume arm-1 is optimal. Now to bound the
 regret we need to find a bound to the #
 of times sub-optimal arm is pulled.

$$\mathbb{E}[n_2] = \mathbb{E}\left[n_2 \mid \begin{array}{c} 0,1 \\ \downarrow \quad \downarrow \\ t=1 \quad t=2 \end{array}\right] P[0,1] + \mathbb{E}\left[n_2 \mid 0,0\right] P[0,0]$$

$$+ \mathbb{E}\left[n_2 \mid 1,0\right] P[1,0] + \mathbb{E}\left[n_2 \mid 1,1\right] P[1,1]$$


 By assuming we only exploit & $t \geq 3$

$$\geq \mathbb{E}\left[n_2 \mid 0,1\right] P[0,1]$$

$$\gg (\tau-2) (1-\mu_1) \mu_2$$

$$\mathbb{E}[n_2] \gg O(\tau)$$

—

- Even if we explore for T_e steps regret for greedy policy will be of order (τ) .

ϵ -GREEDY POLICY:

- Algorithm:
 - With $(1-\epsilon)$ follow greedy policy
 - With (ϵ) prob. choose the arm randomly.

- Regret in this case, $R_T \geq \frac{\epsilon}{2} T$

or
↓

\therefore with atleast ϵ prob.
we take random action - hence
pulling suboptimal arm is of
probability $(\epsilon/2)$.

$$\text{Regret} \gg O(\tau)$$

VARIABLE EPSILON GREEDY:

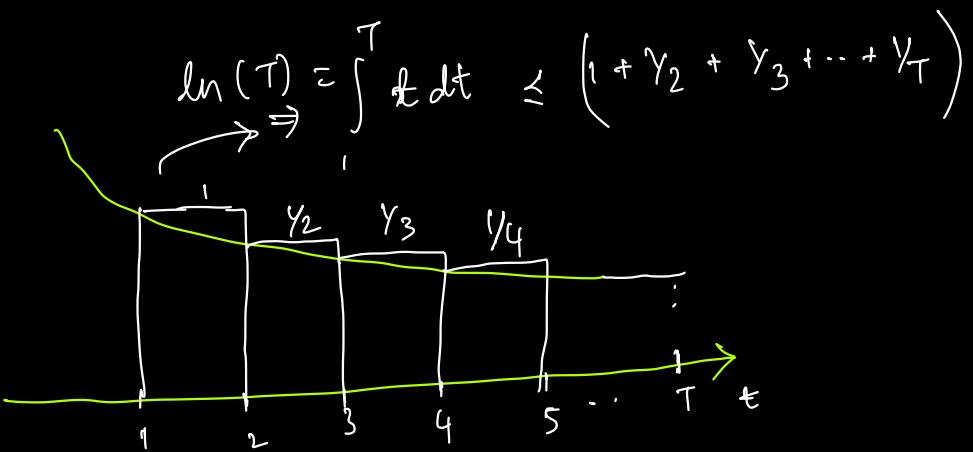
• Markov Chain:

• Instead of fixing ' ϵ ' what time step? if we decrease at each time step? i.e., $\epsilon_t = \frac{\epsilon}{t}$

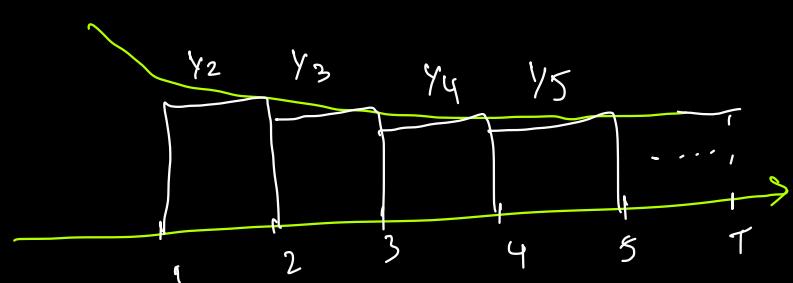
$$\Rightarrow \mathbb{E}[n_2] > \epsilon \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{T} \right)$$

• Bandits:

• Bandit Algorithms
- Tsy and Zaba



• $f(x)$ goes asymptotically as $g(n)$ means the rate at which $f(n)$ changes is same as the rate of change in $g(n)$ as $n \rightarrow \infty$.



$$\ln(T) = \int_1^T t dt \geq \left(\frac{1}{2} + y_3 + \dots + y_T \right)$$

• We can show that,

$$\text{Regret} \leq O(\ln T)$$

Not straight forward

to prove. Need to use concentration inequalities.

We are lower-bounding $\mathbb{E}[n_2]$ with (\cdot) . Then we upperbound (\cdot) by $\ln(T)$. Then how is Regret upper bounded by $O(\ln T)$?

- Refer Auer et al.

SOFTMAX STRATEGY :

We pick each arm with certain probability
(i.e. action isn't deterministic at a time step)

The prob. of picking arm 1 at time step t

$$P\{a_t = 1\} = \frac{\exp\{\hat{\mu}_1(t)\}}{\sum_{i=1}^2 \exp\{\hat{\mu}_i(t)\}}$$

$$P\{a_t = 2\} = \frac{\exp\{\hat{\mu}_2(t)\}}{\sum_{i=1}^2 \exp\{\hat{\mu}_i(t)\}}$$

This has Regret of $O(T)$.

A modification of this is softmax policy
with hyperparameter called as temperature.

$$P\{a_t = 1\} = \frac{\exp\{\hat{\mu}_1(t)/T\}}{\sum_{i=1}^2 \exp\{\hat{\mu}_i(t)/T\}}$$

For high T the action would be uniformly random and for small T the policy would be greedy.

LL/1:

RECAP :

- The expected # of times an arm is pulled till time T is

$$\mathbb{E}[n_i(T)] = \mathbb{E}\left(\sum_{t=1}^T \mathbf{1}_{\{A_t=i\}}\right)$$

$$= \sum_{t=1}^T \mathbb{E}\left[\mathbf{1}_{\{A_t=i\}}\right]$$

$$= \sum_{t=1}^T \mathbb{P}\{A_t=i\}$$

- For ϵ -greedy : $\mathbb{P}\{A_t=i\} \geq \frac{\epsilon}{2}$

- For variable ϵ -greedy : $\mathbb{P}\{A_t=i\} \geq \frac{\epsilon}{t} \xrightarrow{\text{Initial value}}$

$$\Rightarrow \mathbb{E}[n_2(T)] \geq \epsilon \ln T$$

- For K -arms, regret can be written as,

$$\text{Regret} = \mu^* T - \mathbb{E}\left(\sum_{k=1}^K \mu_k \left(\sum_{t=1}^T \mathbf{1}_{\{A_t=k\}}\right)\right)$$

$$= \mu^* T - \sum_{k=1}^K \mu_k \underbrace{\mathbb{E}[n_k]}_{\text{Expected # of times arm } k \text{ played.}}$$

Without loss of generality,

$$\text{Regret} \leq \max_{k \in [K]} (\mu^* - \mu_k) \sum_{k=2}^K E[n_k]$$

Assuming 1st arm is the optimal arm

- Generalizing the algorithms discussed for $[K]$ arms:

GREEDY ALGORITHM :

The estimate of mean reward for arm 'k' at round 't' is,

$$\hat{\mu}_k(t) = \frac{\sum_{z=1}^t R_z(z) \mathbb{1}_{A_z=k}}{\sum_{z=1}^t \mathbb{1}_{A_z=k}}$$

Algorithm : 1) Play each arm once *Explore*

2) Then play the arm with highest empirical estimate. *Exploit*

$$A_{t+1} = \arg \max_{i \in \{1, K\}} \hat{\mu}_i(t)$$

ϵ -GREEDY :

- For $t = 1, 2, \dots, T$,
- $$A_t = \begin{cases} \text{greedy} & \text{Rule w.p } (1-\epsilon) \\ \text{Arm } 'k' & \text{w.p } \frac{\epsilon}{K} \end{cases}$$

VARIABLE ϵ -GREEDY :

- For $t = 1, 2, \dots, T$
- $$A_t = \begin{cases} \text{greedy} & \text{Rule w.p } (1-\epsilon_t) \\ \text{Arm } 'k') & \text{w.p } \frac{\epsilon_t}{K} \end{cases}$$

SOFTMAX :

- For $t = 1, 2, \dots, T$
 - Pick arm ' k' w.p $\hat{\mu}_i(t)/\text{Temp}$
- $$\text{Pr}\{A_t = k\} = \frac{e^{\hat{\mu}_i(t)/\text{Temp}}}{\sum_{j \in \Sigma} e^{\hat{\mu}_j(t)/\text{Temp}}}$$

FOR THE ALGORITHMS :

REGRET BOUNDS

1) Greedy :

$$R \geq O(T)$$

$$\Rightarrow R = O(T)$$

2) ϵ -Greedy :

$$R \geq O(T)$$

$$\Rightarrow R = O(T)$$

3) ϵ_t - Greedy : $R \geq O(\ln T)$
 . We can show that ϵ_t has regret, $R \leq O(\ln T)$
 Refer to Auer. et.al

UPPER CONFIDENCE BOUND (UCB) :

$$P\{ \hat{\mu}_k - \mu > \epsilon \} \leq e^{-\frac{\epsilon^2}{2s^2}}$$

$$A_t = \arg \max_k \left\{ \hat{\mu}_k + b_k \right\}_{k=1}^K$$

↓
Bonus term

$$\text{Here, } b_k(t) = \sqrt{\frac{2 \ln t}{n_k(t)}}$$

$$\overbrace{n_k(t)}$$

$$\overbrace{\ln(t)}$$

. Algorithm :

- 1) Give 1 chance to every arm
- 2) For $t = K+1, \dots, T$ pick the arm
 that has highest UCB $\Rightarrow A_t = \arg \max_k \left\{ \hat{\mu}_k(t) + b_k(t) \right\}_{k=1}^K$

If $\sqrt{\frac{\ln t}{n_k(t)}}$ → very small value as $t \rightarrow \infty$
 Large value and $n_k(t) \rightarrow$ very large value

NOTE :

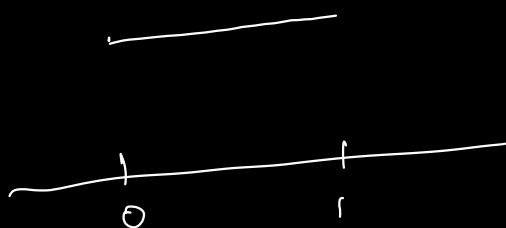
BAYESIAN

ESTIMATION :

Let μ_K be the unknown R.V.
 we assume prior distⁿ for μ_K .
 When the Reward follows the Bernoulli
 distribution we assume the prior distⁿ for μ_K
 to be β -distribution, i.e. $\mu_K \sim \beta(1, 1)$

β - DISTRIBUTION :

$$\beta(1, 1)$$



The PDF of β -dis. is

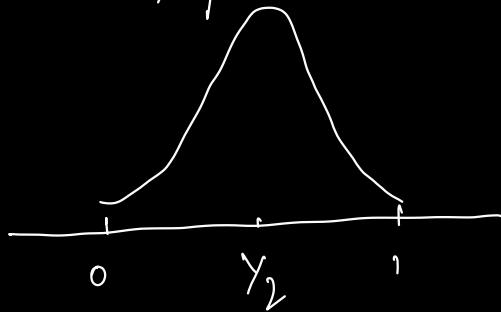
$$f_X(x) = \left\{ \frac{\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}}{x^{\alpha-1}(1-x)^{\beta-1}} \right\} = B(\alpha, \beta) = \int_0^1 x^{\alpha-1}(1-x)^{\beta-1} dx$$

$$\text{where, } \Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$

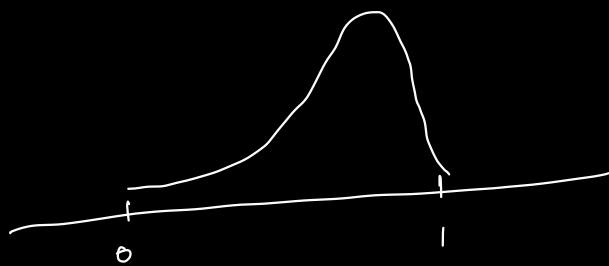
$$\mathbb{E}[X] = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)}$$

$$= \frac{\alpha}{\alpha+\beta}$$

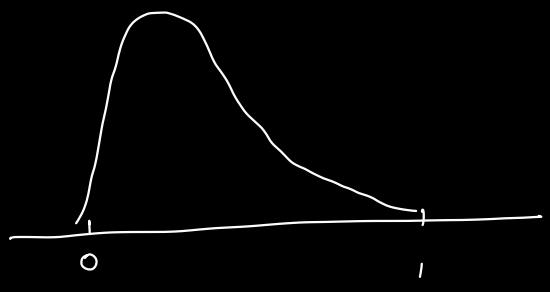
$\alpha = 1000, \beta = 1000$: (Symmetric)



$\alpha = 10,000, \beta = 100$



$\alpha = 100, \beta = 10,000$



1/1:

THOMPSON SAMPLING :

Consider 2-arm case.

Assumption : The mean-reward of each arm is a random variable.

When rewards follow Bernoulli distribution we assume the true rewards come from β -distribution.

$$Hyp : [K] = \{1, 2\}, \mu_1, \mu_2$$

Initialisation : $\mu_1 \sim \beta(1, 1)$, $\mu_2 \sim \beta(1, 1)$

for $t = 1, 2, \dots, T$

sample x_1, x_2 from their respective dist.

$$A_t = \arg \max_i \left\{ x_i \right\}_{i=1}^K$$

Update the dist' of the arm played

$$\text{i.e. } \mu_{A_t} \sim \beta(\alpha_{t-1} + R_t^{A_t}, \beta_{t-1} + 1 - R_t^{A_t})$$

NOTE :

- β - Distribution is a conjugate prior for Bernoulli distribution.

Consider a coin with unknown bias ' p '.

Prior dist. of $p \sim \beta^{(1,1)}$

Assumption:

Given random tosses

outcomes $\{x_1, x_2, \dots, x_n\}$

estimate p .

$$\underbrace{\left\{ \Pr \{D | p=\theta\} \right\}}_{\text{Likelihood}} = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}$$

PDF \rightarrow
not PMF

$$\Rightarrow \underbrace{\Pr \{p=\theta | D\}}_{\text{Posterior dist.}} \propto \Pr \{D | p=\theta\} \Pr \{p=\theta\}$$

$$\propto \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$
$$\propto \theta^{\alpha+\sum x_i-1} (1-\theta)^{\beta+n-\sum x_i-1}$$

$$\Rightarrow \text{PDF } \{p=\theta | D\} \sim \beta(\alpha + \sum x_i, \beta + n - \sum x_i)$$

INVERSE TRANSFORM METHOD :

Let X be the R.V which follows a dist'n. and we need to find a way to sample it from that distribution given that we only use samples from uniform distribution.

CDF of X .

Let $F_X(x)$ be the

Let $U \sim \text{Unif}(0, 1)$

$$F_X(x) = P\{X \leq x\}$$

$$= P\{\tau(U) \leq x\}$$

$$= P\{U \leq \tau^{-1}(x)\}$$

$$F_X(x) = \tau^{-1}(x)$$

$$\Rightarrow \tau^{-1} = F_X$$

$$\Rightarrow \boxed{x = F_X^{-1}(U)}$$

$$f_X(x) = 1 - e^{-x^n} = n$$

$$\tau(u) = x$$

$$\Rightarrow \tau^{-1}(x) = u$$

$$x = F_X^{-1}(u)$$

Let $y \sim \exp(\lambda)$:

$$\therefore F_X(x) = 1 - e^{-\lambda x}$$

$$y = -\frac{1}{\lambda} \ln(1 - v)$$

POLICY BASED BANDIT ALGORITHM :

round 't' we sample arm according

- At each round 't' we sample arm according
- to π_t and update π as, convex combn. of $\pi_t(A_t)$ and 1
- If $R_t = 1$:

$$\pi_{t+1}(A_t) = (1-\alpha) \pi_t(A_t) + \alpha \left\{ \begin{array}{l} \pi_{t+1}(A_t) \\ \pi_t(A_t) \end{array} \right\} + 1$$

$$\pi_{t+1}(a) = (1-\alpha) \pi_t(a) + \alpha (0) \left\{ \begin{array}{l} \pi_{t+1}(a) \\ \pi_t(a) \end{array} \right\} + 1$$

$\pi_{t+1}(a)$ other arm

- If $R_t = 0$:

$$\pi_{t+1}(A_t) = (1-\beta) \pi_t(A_t) + \beta \cdot 0 = \pi_t(A_t) + \beta (0 - \pi_t(A_t))$$

$$\pi_{t+1}(a) = (1-\beta) \pi_t(a) + \beta \cdot 0 = \pi_t(a) + \beta (1 - \pi_t(a))$$

\downarrow
step-size

- When $\alpha = \beta$, the algorithm is called linear reward and penalty (L_{RP}). • When to use L_{R-EP} , L_{R-I} , L_{RP} .
- When $\alpha \gg \beta$, the algorithm is L_{R-EP} .
- When $\beta = 0$, L_{R-I} ; Linear reward inaction.
The update rule for prob. of arm pulled is proportional to reward in linear fashion. (Only when reward is 1 (in-action))

REINFORCE ALGORITHM :

- Here, the prob. of picking an arm at any time t is parameterized by parameter ' θ '.

- Simple parameterization for two-arm case is,

$$\pi_t(1, \theta) = \theta$$

$$\pi_t(2, \theta) = 1 - \theta$$

- Let $p_t(\theta)$ be the average reward obtained by a policy π_θ at time t ,

$$\Rightarrow p(\theta) = \mu_1 \pi(1, \theta) + \mu_2 \pi(2, \theta)$$

$$\begin{aligned}
\Rightarrow \nabla_{\theta} \ell(\theta) &= \mu_1 \nabla_{\theta} \pi^{(1, \theta)} + \mu_2 \nabla_{\theta} \pi^{(2, \theta)} \\
&= \mu_1 \frac{1}{\pi^{(1, \theta)}} \nabla_{\theta} \pi^{(1, \theta)} \pi^{(1, \theta)} + \mu_2 \frac{1}{\pi^{(2, \theta)}} \nabla_{\theta} \pi^{(2, \theta)} \pi^{(2, \theta)} \\
&= \mu_1 \nabla_{\theta} \ln(\pi^{(1, \theta)}) \pi^{(1, \theta)} + \mu_2 \nabla_{\theta} \ln(\pi^{(2, \theta)}) \pi^{(2, \theta)} \\
&= \mu_1 \mathbb{E} \left[\nabla_{\theta} \ln(\pi^{(1, \theta)}) \right] + \mu_2 \mathbb{E} \left[\nabla_{\theta} \ln(\pi^{(2, \theta)}) \right]
\end{aligned}$$

$$\nabla_{\theta} \ell(\theta) = \mathbb{E}_{a \sim \pi_{\theta}} \left[R^{(a)} \nabla_{\theta} \ln(\pi^{(a, \theta)}) \right] = \mathbb{E}_{a \sim \pi_{\theta}} \left[M_a \nabla_{\theta} \ln(\pi^{(a, \theta)}) \right]$$

Update Rule :

$$\theta_{n+1} = \theta_n + \alpha \nabla_{\theta} \ell(\theta_n)$$

The estimate of $\nabla_{\theta} \ell(\theta)$ using the samples is given by,

$$\begin{aligned}
\hat{\nabla}_{\theta} \ell(\theta) &= \frac{1}{n} \sum_{t=1}^n R_t^{(at)} \nabla_{\theta} \ln(\pi^{(at, \theta)}) \\
&= \frac{1}{n} \left[\sum_{i=1}^{n_1} R_i^{(1)} \nabla_{\theta} \ln(\pi^{(1, \theta)}) + \sum_{i=1}^{n_2} R_i^{(2)} \nabla_{\theta} \ln(\pi^{(2, \theta)}) \right] \\
&\quad (n_1 + n_2 = n)
\end{aligned}$$

$$= \sum_{i=1}^{n_1} \frac{R_i^{(1)}}{n_1} \cdot \frac{n_1}{n} \nabla_{\theta} \ln(\pi^{(1, \theta)})$$

$$+ \sum_{i=1}^{n_2} \underbrace{\frac{R_i^{(2)}}{n_2}}_{\xrightarrow{n \rightarrow \infty} \mu_2} \cdot \underbrace{\frac{n_2}{n}}_{\xrightarrow{n \rightarrow \infty} 1} \nabla_\theta \ln(\pi(2, \theta))$$

$$\begin{aligned} \nabla_\theta \hat{P}(\theta) &= \mu_1 \pi(1, \theta) \nabla \ln \pi(1, \theta) \\ &\quad + \mu_2 \pi(2, \theta) \nabla \ln \pi(2, \theta) \end{aligned}$$

\approx

stochastic gradient ascent :

The stochastic gradient ascent :

$$\theta_{n+1} = \theta_n + \alpha R_i^{(a)} \nabla \ln \pi(a, \theta)$$

REINFORCE ALGORITHM WITH BASELINE :

The update rule is,

$$\theta_{n+1} = \theta_n + \alpha (R^{(a)} - b_n) \nabla \ln \pi(a, \theta)$$

where, $b_n = \frac{1}{n} \sum_{i=1}^n R_i$

$$\theta_K = \ln \left[\sum \exp \{ \theta_j \} \right]$$

$$\exp \{ \theta_K \} \\ () \\ \pi_K$$

2Y₁:

REINFORCE ALGORITHM FOR K-ARMS :

- For arm k ,

$$\pi(k, \theta) = \frac{\exp\{\theta_k\}}{\sum_{j=1}^k \exp\{\theta_j\}}$$

- Update Rule :

$$\theta_{n+1} = \theta_n + \nabla P(\theta)$$

CONCENTRATION INEQUALITY :

MARKOV INEQUALITY :

For a non-negative R.V, X

$$P\{X > a\} \leq \frac{E X}{a}$$

CHEBYSHEV'S INEQUALITY :

for a R.V with finite mean and variance,

$$P\{|X - \mu| > \varepsilon\} \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

WLLN :

A dist. is heavy

- Let X_1, X_2, \dots, X_n be i.i.d R.V's tail if $\underbrace{\text{M.g.F}}_{\mathbb{E}[e^{\lambda X}]} = \infty$

then

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| > \varepsilon\right\} \xrightarrow{n \rightarrow \infty} 0$$

(if λ)

Exponen-

tial dist. is

heavy tail

Proof:

From Chebychev's inequality,

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| > \varepsilon\right\} \leq \frac{\sigma^2}{n\varepsilon^2}$$

CHEBNOFF

HOEFFDING

INEQUALITY:

- Let X_1, X_2, \dots, X_n be i.i.d R.V with mean μ .

Then,

$$\mathbb{P}\{|S_n - \mu| > \varepsilon\} \leq 2 \exp\{-n\varepsilon^2\}$$

UCB ALGORITHM :

- For $t=1, \dots, K$ play a_K at round 'K'

- For $t = K, \dots, T$

$$A_t = \arg \max_K \hat{\mu}_{K,t-1} + \sqrt{\frac{2 \log t}{n_{K,t-1}}}$$

$\underbrace{\quad}_{\text{UCB}_{K,t-1}}$

• Claim: $R_T(\text{UCB}) \leq O(\log T)$

Proof:

$$R_T(\text{UCB}) = \mathbb{E} \left[\sum_{t=1}^T \mu^* - \mu_{A_t} \right]$$

↓
Arm played at round t

$$= \sum_{k=1}^K \mathbb{E}[n_k] \Delta_k$$

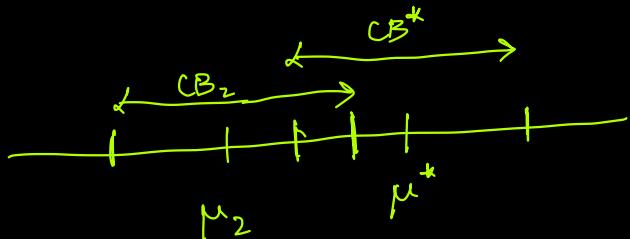
~~~~~  
Expected # of times we play  
arm ' $k$ '

• Let  $n_i(T)$  be the # of times arm ' $i$ ' is picked.

$$\Rightarrow n_i(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t = i\}}$$

• At time ' $t+1$ ' if we play arm ' $k$ '

$$\hat{\mu}_t^* + \sqrt{\frac{2 \ln t}{n_{*,t}}} < \hat{\mu}_k + \sqrt{\frac{2 \ln t}{n_{k,t}}}$$



For the above inequality to hold true one of the following cases should be true:

Case 1:  $\hat{\mu}_* < \mu_* - \sqrt{\frac{2 \ln t}{n_{*,t}}} \quad (\text{Underestimate optimal arm})$

Case 2:  $\hat{\mu}_K > \mu_K + \sqrt{\frac{2 \ln t}{n_{K,t}}} \quad (\text{Overestimate of sub-optimal arm})$

Case 3:  $n_i < \frac{8 \ln T}{\Delta_i^2} \quad (\text{less exploration of arm } i)$

Claim:

Suppose none of the above cases are true, then  $i$ th arm will n't get picked at time  $t+1$ .

Proof:  $\hat{\mu}_* > \mu_* - \sqrt{\frac{2 \ln t}{n_{*,t}}} \quad (\text{case 1 not true})$

$$\begin{aligned} \hat{\mu}_* + \sqrt{\frac{2 \ln t}{n_{*,t}}} &> \mu_* \\ &> \mu_* - \mu_i + \mu_i \\ &> \Delta_i + \mu_i \end{aligned}$$

$$\geq \sqrt{\frac{8 \ln t}{n_{i,t}}} + \mu_i \quad (\text{Case 3 not done})$$

$$\geq \hat{\mu}_i + \sqrt{\frac{2 \ln t}{n_{i,t}}} + \sqrt{\frac{2 \ln t}{n_{i,t}}}$$

$$\hat{\mu}_* + \sqrt{\frac{2 \ln t}{n_{*,t}}} \geq \hat{\mu}_i + \sqrt{\frac{2 \ln t}{n_{i,t}}} \quad (\text{Case 2 not done})$$

$\underbrace{\phantom{\hat{\mu}_* + \sqrt{\frac{2 \ln t}{n_{*,t}}}}}_{UCB_{*,t}}$

$$\geq UCB_{i,t}$$

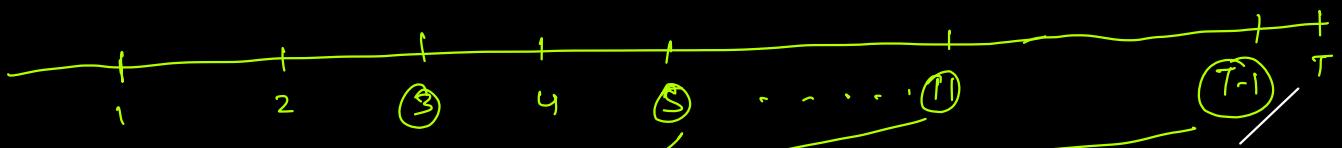
Now consider,

$$\mathbb{E}[n_i(T)] = \mathbb{E}\left[\sum_{i=1}^n \mathbb{I}\{A_t=i\}\right]$$

splitting the action taken at each time on a event,

$$\left\{ \begin{array}{l} \mathbb{E}\left[\sum_{i=1}^T \mathbb{I}\{A_t=i, n_i(t) < u(t)\}\right] \\ + \mathbb{E}\left[\sum_{i=1}^T \mathbb{I}\{A_t=i, n_i(t) \geq u(t)\}\right] \end{array} \right.$$

$u(t) = \frac{8 \ln t}{\Delta_i^2}$



At these times  $\{n_i(t) \geq u(t)\}$   
and remaining times  $\{n_i(t) < u(t)\}$

$$\text{WLT} \quad \text{if} \quad n_i(t) < \frac{8 \ln T}{\Delta_i^2} = u(T)$$

$$\text{then } ① \leq \frac{8 \ln T}{\Delta_i^2}$$

Also  $\forall t$  such that  $n_i(t) > \frac{8 \ln T}{\Delta_i^2}$

$$n_i(t) > \frac{8 \ln T}{\Delta_i^2} > \frac{8 \ln t}{\Delta_i^2}$$

$$\Rightarrow E[n_i(T)] \leq \frac{8 \ln T}{\Delta_i^2} + \sum_{t=1}^T \underbrace{\mathbb{P}\{A_t \subseteq i, n_i(t) > u(t)\}}_{\substack{\text{only possible either due} \\ \text{to case 1 or 2.}}}$$

$$\leq \frac{8 \ln T}{\Delta_i^2} + \sum_{t=1}^T \underbrace{\mathbb{P}\{\text{case 1 or case 2 occurs}\}}_{\substack{}}$$

$$\leq \frac{8 \ln T}{\Delta_i^2} + \sum_{t=1}^T \underbrace{\mathbb{P}\{\text{case 1 occurs}\}}_{\substack{}}$$

$$+ \sum_{t=1}^T \underbrace{\mathbb{P}\{\text{case 2 occurs}\}}_{\substack{}}$$

From Hoeffding's inequality,

$$\mathbb{P}\{\text{case 1 is true}\} = \mathbb{P}\{\hat{\mu}_* \leq \mu_* - \sqrt{\frac{2\ln t}{n_{*,t}}}\}$$

$$\leq \exp\left\{-n_{*,t} \cdot \frac{2\ln t}{n_{*,t}}\right\}$$

$$= \gamma_t^2$$

$$\text{W.H., } \mathbb{P}\{\text{case 2 is true}\} = \mathbb{P}\{\hat{\mu}_i > \mu_i + \sqrt{\frac{2\ln t}{n_{i,t}}}\}$$

$$\leq \exp\left\{-n_{i,t} \cdot \frac{2\ln t}{n_{i,t}}\right\}$$

$$= \gamma_t^2$$

$$\Rightarrow \mathbb{E}[n_i(T)] \leq \frac{8\ln T}{\Delta_i^2} + \sum_{t=1}^T \gamma_t^2$$

$$\leq \frac{8\ln T}{\Delta_i^2} + \underbrace{(1 - \gamma_T)}_c$$

$$\Rightarrow \mathbb{E}[R_T(\text{UCB})] \leq O(\ln T)$$

NOTE :

- There is a lower bound on the regret.
- For MAB (Stochastic) no algorithm can incur a regret less than logarithmic regret. . Lai - Robinson

$$\text{i.e. } \mathbb{E}[n_i(T)] \geq O(\ln T)$$

PAC OPTIMALITY :

- For a PAC( $\epsilon, \delta$ ) ,

Input :  $\epsilon, \delta$

Output : One arm  $j$  such that,

$$P\{\mu_j > \mu^* - \epsilon\} \geq 1 - \delta$$

- In other words, probability of not picking an  $\epsilon$ -close arm is  $\leq \delta$ .

# NAIVE ROUND ROBIN ALGORITHM :

Algorithm :

- a certain no. of rounds.
- Play all arms in a round robin fashion for a certain no. of rounds.
- Declare the arm with highest empirical estimate as the optimal arm.

NOTE :

- Sample Complexity :  
How many rounds should we play each arm such that the arm we output is  $\epsilon$ -close to the optimal arm with a high-probability.

Claim :

- Naive - Round Robin algorithm is PAC optimal.



Let  $N = \frac{K}{\varepsilon^2} \log\left(\frac{K}{\delta}\right)$ . Then the arm that we output will be  $\varepsilon$ -close to the optimal arm with a very high probability.

Also, each arm is played  $\frac{1}{2\varepsilon^2} \log\left(\frac{2K}{\delta}\right)$  times.

Let  $\{\hat{\mu}_k\}_{k=1}^K$  be the estimate of rewards after 'N' rounds. Then choose arm  $j = \arg \max_{k \in [K]} \hat{\mu}_k$

W.R.T, Picking arm ' $j$ ' will be  $\underbrace{\text{bad if}}_{\text{In relation with PAC bound.}}$ ,

$$\mu_j < \mu^* - \varepsilon$$

We want to bound,

$$\Pr\{\mu_j < \mu^* - \varepsilon\}$$

Let 'l' be the arm such that,  $\mu_l < \mu^* - \varepsilon$ .

$\therefore$  Probability of picking the  $l^{\text{th}}$  arm at the end

of 'N' rounds is,

$$\Pr\{\hat{\mu}_l \geq \{\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K\}\}$$

28/1:

Quiz 1 : Feb 4

Recap :

Assignment 2 : Feb 2  
(11:30 AM)

In PAC optimality we are interested in finding the no. of trials required to discover an  $\epsilon$ -optimal arm with high probability of  $(1-\delta)$ .

i.e. determine  $N$  such that ;

$$\Pr \{ \mu_j < \mu_* - \epsilon/2 \} \leq \delta$$

$\downarrow$   
True reward of returned 'j'.

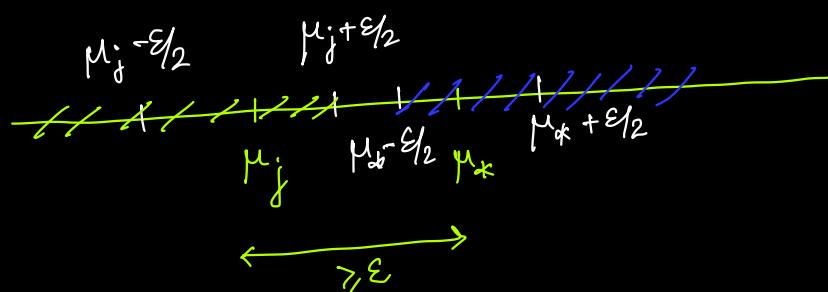
arm 'j' if ,

We pick

$$\Pr \{ \hat{\mu}_j > \hat{\mu}_* \} \leq \Pr \{ \hat{\mu}_j > \mu_j + \epsilon/2 \text{ or } \hat{\mu}_* \leq \mu_* - \epsilon/2 \}$$

Explanation for this :

Suppose not : i.e.  $\hat{\mu}_j < \mu_j + \epsilon/2$  and  $\hat{\mu}_* > \mu_* - \epsilon/2$



$$\Rightarrow \hat{\mu}_* > \hat{\mu}_j$$

$$\cdot \Rightarrow P\{\hat{\mu}_j > \hat{\mu}_*\} = P\{\hat{\mu}_j > \mu_j + \varepsilon_2\} \cup P\{\hat{\mu}_* \leq \mu_* - \varepsilon_2\}$$

$$\Rightarrow P\{\hat{\mu}_j > \hat{\mu}_*\} \leq \underbrace{P\{\hat{\mu}_j > \mu_* + \varepsilon_2\}}_{\text{Apply Hoeffding's inequality}} + P\{\hat{\mu}_* \leq \mu_* - \varepsilon_2\}$$

$$\leq \exp\{-2N_j\varepsilon^2\} + \exp\{-2N_*\varepsilon^2\}$$

$$\cdot \text{ w.r.t., } N_j = N_* = \frac{1}{2\varepsilon^2} \ln\left(\frac{2K}{\delta}\right)$$

$$\Rightarrow P\{\hat{\mu}_j > \hat{\mu}_*\} \leq 2 \exp\{-2 \ln\left(\frac{K}{\delta}\right)\}$$

$$\leq \frac{\delta}{K}$$

$$\Rightarrow \text{Total probability of picking the arm 'i' with } \mu_i < \mu_* - \varepsilon$$

$$\text{is, } P\left\{\bigcup_{i \neq *} \{\hat{\mu}_i > \hat{\mu}_*\}\right\} \leq \sum_{i \neq *} \underbrace{P\{\hat{\mu}_i > \hat{\mu}_*\}}_{\delta/K} \leq \frac{K \cdot \delta}{K} = \underline{\delta}$$

$$\leq \frac{K \cdot \delta}{K} = \underline{\delta}$$

∴ Naive algorithm is PAC optimal with ,

$$N = \frac{K}{2\epsilon^2} \ln\left(\frac{2K}{\delta}\right)$$

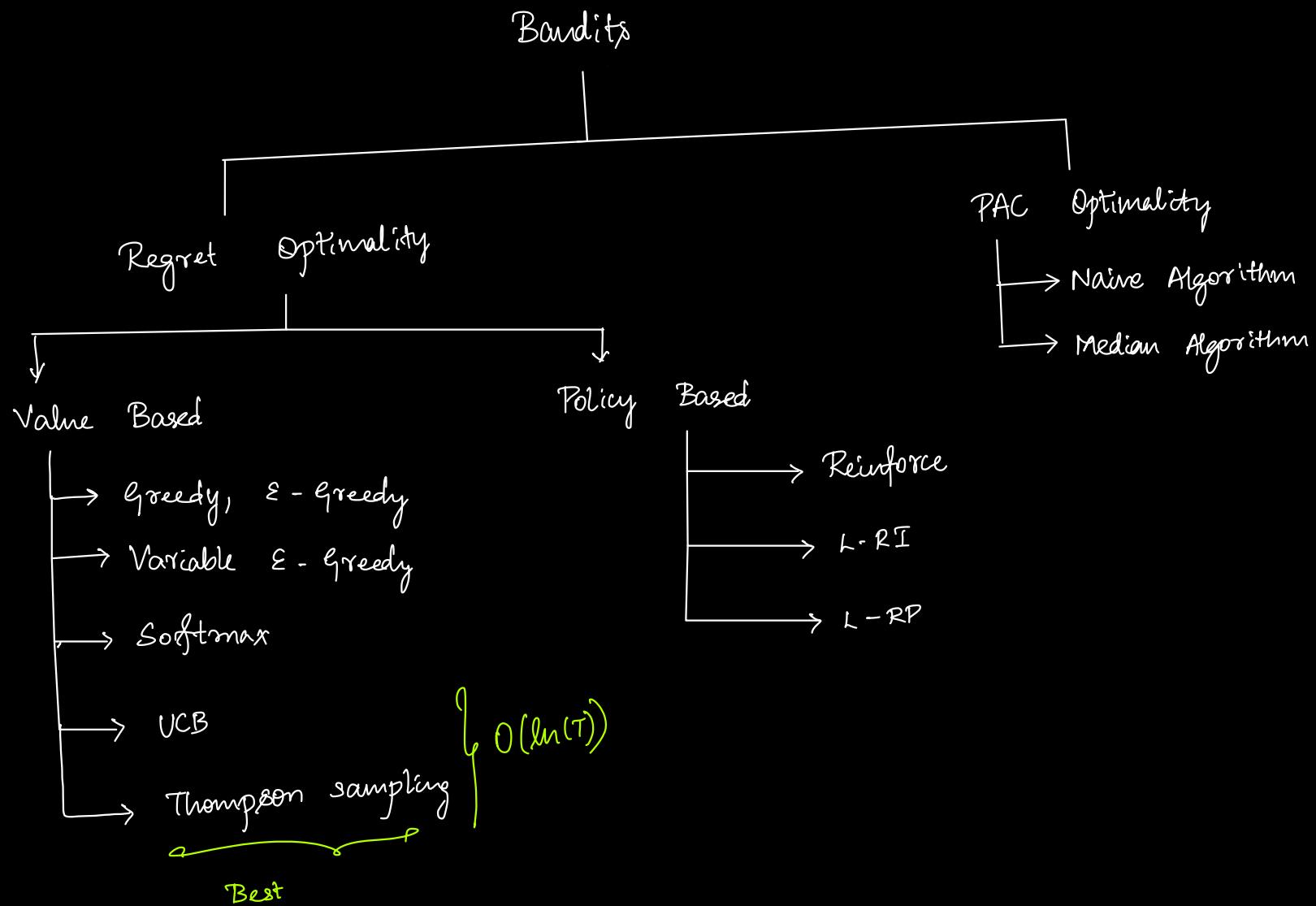
NOTE :

Median elimination algorithm has better PAC optimality

$$\therefore N_{\text{median}} \leq N_{\text{Naive}}$$

## BANDITS

## SUMMARY :



31/1 :

MARKOV DECISION PROCESS : (think this as controlled Markov chain)

• 4 - Tuple  $(S, A, P, R)$

•  $P : S \times S \times A \rightarrow [0, 1]$

where  $P_{ij}(a) :$  Probability of starting from 'i' and ending up in 'j' when action 'a' is taken.

•  $R : S \times A \times S \rightarrow \mathbb{R}$

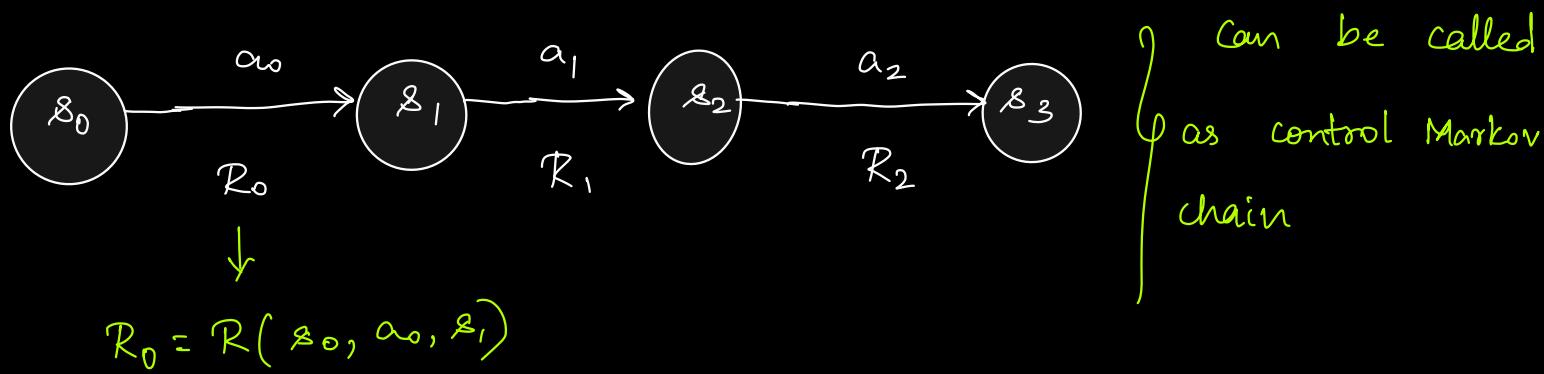
where,  $R_{ij}(a) :$  Reward you will get starting from state 'i' and transitioning to 'j' when choosing action 'a'.

If the states and actions are finite then we

say MDP is finite MDP.

• For each action  $a$ ,  $P(a)$  is a  $|S| \times |S|$  matrix

Segmental decision and reward process:



We want MDP to achieve maximum long-term

reward.

TOTAL REWARD CRITERION:

We want to maximize,  $\mathbb{E} \left[ \sum_{i=1}^T R_i \right]$

- Here we need a special state ' $s_T$ ' in the MDP (terminal state) which is an absorbing state

$$\text{with } R(s_T, \cdot, s_T) = 0$$

CRITERION:

TOTAL DISCOUNTED REWARD CRITERION:

Let  $\gamma \in [0, 1)$  be the discount factor.

$$\text{Goal : } \max \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R_k \right]$$

Using discounted reward results in :

- 1) Finite return
  - 2) No need to have a terminal state
  - 3) Weightage for a farthest reward to decreases.
- (i.e. immediate emphasis)*
- rewards are given more

AVERAGE REWARD CRITERION :

Goal :  $\max_{T \rightarrow 0} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^T R_t \right]$

Example :

|   | 1 | 2 | 3  | 4    |
|---|---|---|----|------|
| 1 | 0 | 0 | 0  | 1    |
| 2 | 0 | 6 | 0  | -100 |
| 3 | 0 | 0 | 0  | 0    |
| 4 | 8 | 9 | 10 | 11   |

Consider,

Let  $\gamma = 0.9$

- State space : Cells and  $|S| = 11$
- Action space : { Left, Right, Up, Down }
- Reward :

  - $R(i, \cdot, \cdot) = 0 \quad \forall i \neq 4, 7$
  - $R(4, \cdot, \cdot) = 1$
  - $R(7, \cdot, \cdot) = -100$

- Probability Transition Matrix :
  - If there is a wall then you remain in the same state.
  - 80% of time you are going to the intended state for a given action.
  - 10% of time you end up in a cell in the perpendicular direction.

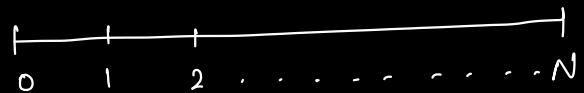
If we define a terminal state for the above example then we can consider total reward instead of discounted reward.

### INFINITE HORIZON PROBLEM :



Here usually rewards are independant of time 't'.

### FINITE HORIZON PROBLEM :



Here it's possible to formulate the reward as a function of time.

The Bandit setting has no 'Markov nature'  
(No notion of states, P(M))

# FINITE DISCOUNTED REWARD MDP :

- Representation :  $\langle S, A, P, R, \gamma \rangle$

- goal :  $\max \mathbb{E} [R_0 + \gamma R_1 + \gamma^2 R_2 + \dots]$

such that expected discounted

- Find  $a_0, a_1, a_2, \dots$

reward is maximized.

NOTION OF A POLICY :

$\leftarrow$  joint strategy for all the actions  
deterministic policy  $\pi : S \rightarrow A$ .

- Stationary

- How many stationary

MDP ?

$$\text{sol}: - |A|^{|S|}$$

deterministic policies in a finite

2/2:

Recap:

Goal in the MDP setting:

- Find an action sequence  $\{a_0, a_1, a_2, \dots\}$  such that it maximizes some long run reward criterion.

$$R_t = R(s_t, a_t, s_{t+1})$$

Notion of policy:

1) stationary

2) Deterministic

Given a policy  $\mu$ , we associate a performance metric of  $\mu$  as,

$$V_\mu(s_0) = \mathbb{E} \left[ R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \mid \underbrace{s_0}_{\text{starting state}} \right]$$

where,  $R_t = R(s_t, \mu(s_t), s_{t+1})$

We define,  $V_\mu = [V_\mu^{(1)}, V_\mu^{(2)}, \dots, V_\mu^{(N)}]$

Here,  $V_\mu$  has a POSET relation.

NOTE :

• POSET Relation :

- Any two elements in the set are comparable.  
(with  $\leq$ )

defn: 1)  $a \leq a$  : Reflexivity

2) If  $a \leq b$ ,  $b \leq a \Rightarrow a = b$

3) If  $a \leq b$ ,  $b \leq c \Rightarrow a \leq c$

- For a POSET relation, there exists

a  $\mu^*$  such that,

$$V_{\mu^*}(i) > V_{\mu}(i) \quad \forall \mu, \forall i \in \{1, \dots, N\}$$

If we can define either  $a \leq b$  or  $b \leq a$   $\forall a, b \in$  set then it is totally ordered.

- Let  $\Pi$  be the set of all deterministic policies.

Then,

$$\mu^*(i) = \arg \max_{\mu \in \Pi} V_{\mu}(i)$$

deterministic policy exist?

- Why should a optimal stationary deterministic policy exist?
  - TPM is stationary / homogeneous  
Is this usage correct?
  - Reward function is stationary

: of these assumption we will have stationary deterministic policy

NOTE :

- Homogeneous Markov Chain :

$$\Rightarrow \mathbb{P}\{X_n = i \mid X_{n-1} = j\} = \mathbb{P}\{X_{n+k} = i \mid X_{n-1+k} = j\}$$

- For a policy  $\mu$ :

$$V_\mu(s_0) = \underbrace{\mathbb{E}[R(s_0, \mu(s_0), s_1)]}_{\substack{\text{w.r.t randomness in next} \\ \text{state and action} \\ \text{if randomized}}} + \gamma \mathbb{E}[R(s_1, \mu(s_1), s_2)] + \gamma^2 \mathbb{E}[R(s_2, \mu(s_2), s_3)] + \dots$$

$$\text{WKT, } \mathbb{E}[R(s_0, \mu(s_0), s_1)] = \sum_{s_1=1}^N P_{s_0, s_1}(\mu(s_0)) R(s_0, \mu(s_0), s_1)$$

Expected immediate reward for policy ' $\mu$ '.

- We can observe that,

$$V_\mu(s_0) = \underbrace{\mathbb{E}[R(s_0, \mu(s_0), s_1)]}_{\substack{\text{Expected Immediate} \\ \text{reward following} \\ \text{policy ' $\mu$ '.}}} + \gamma \underbrace{\sum_{s_1=1}^N P_{s_0, s_1}(\mu(s_0)) V_\mu(s_1)}_{\substack{\text{Expected future reward} \\ \downarrow \\ \mathbb{E}[R(s_1, \mu(s_1), s_2)] + \gamma \mathbb{E}[R(s_2, \mu(s_2), s_3)] + \gamma^2 \mathbb{E}[R(s_3, \mu(s_3), s_4)] + \dots}}$$

The above is true since ,

$$\begin{aligned}
 V_\mu(s_0) &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E} \left[ R(s_t, \mu(s_t), s_{t+1}) \right] \\
 &= \mathbb{E} \left[ R(s_0, \mu(s_0), s_1) \right] + \sum_{t=1}^{\infty} \gamma^t \mathbb{E} \left[ R(s_t, \mu(s_t), s_{t+1}) \right] \\
 &= \mathbb{E} \left[ R(s_0, \mu(s_0), s_1) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E} \left[ R(s_t, \mu(s_t), s_{t+1}) \right] \right] \\
 &= \mathbb{E} \left[ R(s_0, \mu(s_0), s_1) + \gamma \sum_{t=0}^{\infty} \gamma^t \mathbb{E} \left[ R(s_{t+1}, \mu(s_{t+1}), s_{t+2}) \mid s_{t+1} = s' \right] \right] \\
 &\quad \text{conditional expectation on next state} \\
 &\quad \underbrace{V_\mu(s')}_{\text{Value function}}
 \end{aligned}$$

$$\Rightarrow V_\mu(s_0) = \mathbb{E} \left[ R(s_0, \mu(s_0), s_1) \right] + \gamma \sum_{s' \in S} P_{s_0, s'}(\mu(s_0)) V_\mu(s')$$

$$\begin{aligned}
 &= \sum_{s' \in S} P_{s_0, s'}(\mu(s_0)) R(s_0, \mu(s_0), s') \\
 &\quad + \gamma \sum_{s' \in S} P_{s_0, s'}(\mu(s_0)) V_\mu(s')
 \end{aligned}$$

} Possible to write  
 this : PTM is  
 homogeneous and also  
 reward function if  
 stationary.

$$\text{WKT, } V_{\mu}(s_0) = \mathbb{E}[R_0] + \gamma \mathbb{E}[R_1] + \gamma^2 \mathbb{E}[R_2] + \dots$$

$$\mathbb{E}[R_1] = \mathbb{E}[R(s_1, \mu(s_1), s_2)]$$

$$= \mathbb{E}_{s_1} \left[ \mathbb{E}_{s_2 | s_1=j} [ R(s_1, \mu(s_1), s_2) | s_1=j ] \right] = \sum_{s_1} P_{s_0, s_1}(\mu(s_0)) \cdot (\cdot)$$

$$\text{Why } \mathbb{E}[R_2] = \sum_{s_2} P_{s_0, s_2}^2(\mu(s_0)) \mathbb{E}_{s_3 | s_2} [ R(s_2, \mu(s_2), s_3) | s_2 ]$$

$$\text{Let, } R(s, \mu(s)) = \sum_{s' \in S} P_{s, s'}(\mu(s)) R(s, \mu(s), s') = R_{\mu(s)}$$

$$\begin{bmatrix} V_{\mu^{(1)}} \\ V_{\mu^{(2)}} \\ \vdots \\ V_{\mu^{(N)}} \end{bmatrix} = \begin{bmatrix} R_{\mu^{(1)}} \\ R_{\mu^{(2)}} \\ \vdots \\ R_{\mu^{(N)}} \end{bmatrix} + \underbrace{\gamma \cdot P_{\mu}}_{S \times S} \begin{bmatrix} V_{\mu^{(1)}} \\ V_{\mu^{(2)}} \\ \vdots \\ V_{\mu^{(N)}} \end{bmatrix}$$

$$P_{\mu} = P_{i,j}(\mu^{(i)}) = \begin{bmatrix} P_{1,1}(\mu^{(1)}) & P_{1,2}(\mu^{(1)}) & \dots & \dots \\ P_{2,1}(\mu^{(2)}) & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$$\Rightarrow V_\mu = R_\mu + \gamma P_\mu V_\mu$$

$$\Rightarrow V_\mu = (I - \gamma P_\mu)^{-1} R_\mu \quad \begin{array}{l} \text{Policy evaluation} \\ \text{for policy } \mu \\ \downarrow \text{map to } \left( \frac{1}{1-\gamma} \right) \\ (I + \gamma P_\mu + \gamma^2 P_\mu^2 + \gamma^3 P_\mu^3 + \dots) R_\mu \\ = R_\mu + \gamma P_\mu R_\mu + \gamma^2 P_\mu^2 R_\mu + \dots \end{array}$$

why should  $(I - \gamma P_\mu)$  be invertible?

If we show the eigen values of  $P_\mu$  is  $\leq 1$   
 then it is obvious that eigen values of  $(I - \gamma P_\mu) > 0$

Let  $\lambda$  be some eigen value of  $P_\mu$ .

$$\Rightarrow P_\mu \cdot \alpha = \lambda \alpha$$

$$\text{let } x_n = \frac{\alpha}{\|\alpha\|}$$

$$\Rightarrow P_\mu x_n = \lambda x_n$$

Consider  $|x_j| = \max \{|x_1|, \dots, |x_n|\}$

$$\Rightarrow \underbrace{P_\mu(j)}_{1 \times n \text{ vector}} x_n = \lambda x_j$$

$$\begin{aligned}
 \Rightarrow |\lambda x_j| &= \left| \sum_i P_\mu(j, i) x_i \right| \\
 &\leq \sum_i P_\mu(j, i) |x_i| \\
 &\leq \underbrace{\sum_i P_\mu(j, i)}_1 x_j \\
 &\leq |x_j|
 \end{aligned}$$

$$\Rightarrow |\lambda| \leq 1$$

$\Rightarrow$  Max. eigen value is 1.

How to find the optimal policy?

- Naive algorithm is to find  $v_\mu$  w/  $\mu \in \overline{\mathcal{P}}$
- choose  $v_{\mu^*}$  by comparison

NOTE :

- For a matrix  $A \in \mathbb{C}^{n \times n}$  with entries  $a_{ij}$ . Let  $R_i$  be the sum of elements in the  $i^{\text{th}}$  row. Then, every eigenvalue of  $A$  lies within at least one of the  $n$  discs  $D(a_{ii}, R_i)$ .

$$D(a_{ii}, R_i)$$

$\downarrow$   
 $\underbrace{\quad}_{\text{Centre}}$        $\underbrace{\quad}_{\text{Radius}}$

a/2 :

Recap :

- MDP :  $\leftarrow$  Triple  $\langle S, A, P, R \rangle$

formulation :

goal      {  
    | discounted reward       $\max_{a_0, a_1, a_2, \dots} \mathbb{E} [R_0 + \gamma R_1 + \gamma^2 R_2 + \dots]$

deterministic policy :

- $\mu : S \rightarrow A$   
 $\underbrace{|A|}_{(S)}$  different policies

We associate value function,  $V_\mu$  to each policy  $\mu$ .

$$\text{WKT, } V_\mu^{(i)} = \mathbb{E}_\mu \left[ R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \right]$$

$\downarrow$

$$R(i, \mu(i), s_i)$$

$$V_\mu^{(i)} = \begin{matrix} \text{Expected Immediate} \\ \text{Reward} \end{matrix} + \begin{matrix} \text{Expected} \\ \text{Future reward} \end{matrix}$$

Here, Expected Immediate Reward :

$$R_{\mu^{(i)}} = \sum_{j=1}^n P_{i,j}(\mu^{(i)}) R(i, \mu^{(i)}, j)$$

$$\begin{matrix} \text{Expected} & \text{Future} & \text{Reward} \end{matrix}$$

$$= \gamma \sum_{j=1}^n P_{i,j}(\mu^{(i)}) V_{\mu^{(j)}}$$

$$\therefore V_\mu = R_\mu + \gamma P_\mu \cdot V_\mu$$

Conditions for  
stationary deterministic  
policy to be optimal.

- Interpretation of  $V^*$  :
  - $V^*(i)$  : Value function for the optimal policy  $\mu^*$  in state 'i'
  - $V^*(i)$  : Maximum total discounted reward we can collect starting from state 'i'.
- WKT,  $\underbrace{\text{Not a linear eqn.}}_{\text{eqn.}} \quad \left\{ \begin{array}{l} V^*(i) = \max_a \left[ \underbrace{\mathbb{E}[R(i, a, \cdot)]}_{= R(i, a)} + \gamma \sum_{j=1}^n P_{ij}(a) V^*(j) \right] \\ \mu^*(i) = \operatorname{argmax}_a \left[ R(i, a) + \gamma \sum_{j=1}^n P_{ij}(a) V^*(j) \right] \end{array} \right.$
- $\mu^*$  : Optimal stationary deterministic policy.

### BELLMAN OPERATOR :

Bellman Operator ,  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$

Also,  $TV^* = V^*$

Here,  $T(J(i)) = \max_a \left[ R(i, a) + \gamma \sum_{j=1}^n P_{ij}(a) J(j) \right]; \forall i \in S$

We can show that  $v^*$  is unique.

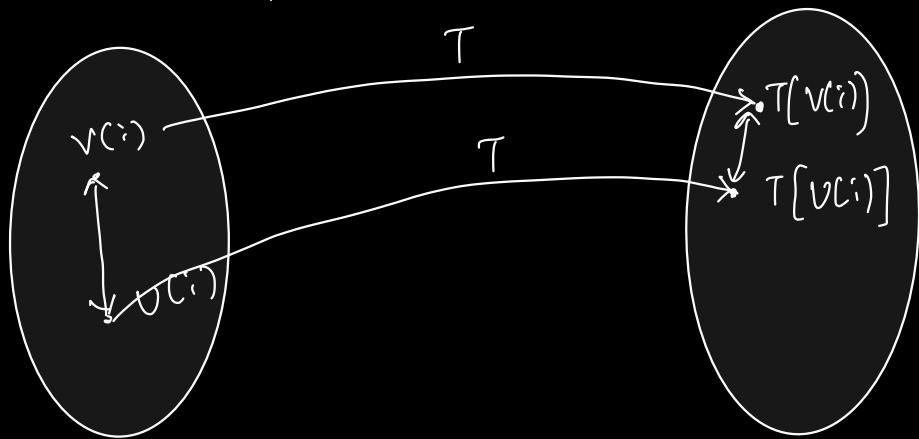
$$\therefore e \quad T(v) \neq v \quad \forall v \neq v^*$$

$v^*$  is called the fixed point of operator 'T'.

The Bellman operator,  $T$  has a special property

called contraction. i.e.,

$$|T(v(i)) - T(v(i'))| \leq \alpha |v(i) - v(i')|$$



The iterative procedure of applying 'T' can be done as a fixed point iteration.

WKT,

$$\left. \begin{array}{l} v_1 = T(v_0) \\ v_2 = T(v_1) \\ \vdots \\ v_n = \underbrace{T^n(v_0)}_{T(\tau(\dots))} \end{array} \right\} \Rightarrow \lim_{n \rightarrow \infty} T^n(v_0) = v^*$$

We need to show that 'T' has contraction property.

## VALUE ITERATION :

Initialization : Initialize  $v_0$

For  $t = 1, 2, \dots$

$$v_t = T(v_{t-1}) \\ = \max_a \left( R(i, a) + \gamma \sum_{j \in S} P_{ij}(a) v_{t-1}(j) \right) \quad \forall i$$

Perform the above step until,  $T(v) = v = v^*$   
In proactive check  $|v_t - v_{t-1}| < \epsilon$  then  
 $v^* = v_t$

W1:

~~1) a<sub>1</sub>/2~~

RECAP:

~~2) w<sub>2</sub>~~

• Bellman Equation :

~~3) w<sub>2</sub>~~

$$V^*(i) = \max_a \left[ R(i, a) + \sum_{j=1}^n P_{ij}(a) V^*(j) \right]$$

~~4) b<sub>2</sub>~~

~~5) b<sub>2</sub>~~

• Optimal Value function :  $V^*$

$$R(i, a) = \sum_{j \in S} P_{ij}(a) R(i, a, j)$$

• Bellman Operator ,  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$   
~~Contraction operator~~

$$T[V(i)] = \max_a \left[ R(i, a) + \sum_{j=1}^n P_{ij}(a) V(j) \right]$$

•  $T(V^*) = V^*$  and  $V^*$  is unique.

STATE ACTION VALUE FUNCTION :

$$Q_\mu : S \times A \rightarrow \mathbb{R}$$

$\forall i \in S, a \in A$

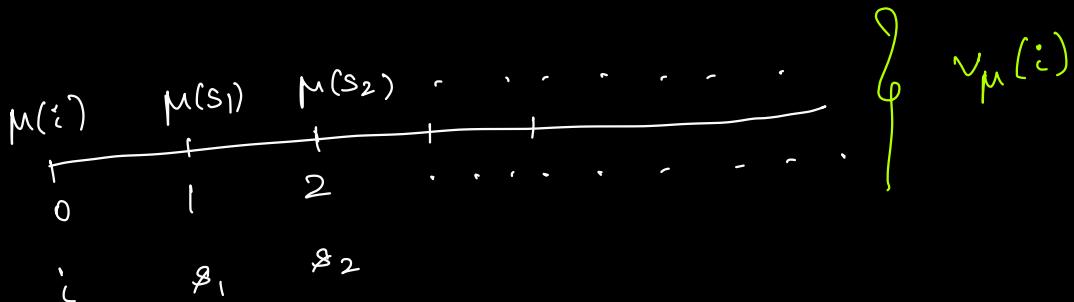
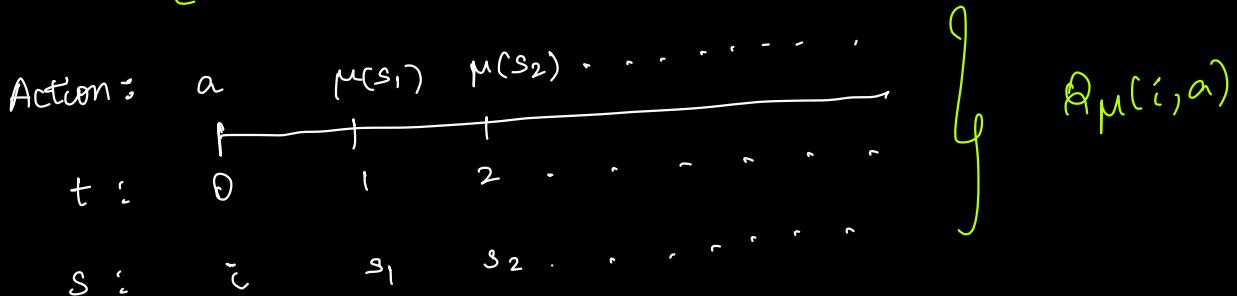
$$\Rightarrow Q_\mu(i, a) \in \mathbb{R}$$

$Q_{\mu}(i, a)$  is the expected reward we obtain by taking action 'a' in state 'i' and following policy  $\mu$  from the next state.

$$Q_{\mu}(i, a) = \mathbb{E}[R(i, a, \cdot)] + \gamma \sum_{j \in S} P_{i,j}(a) V_{\mu}(j)$$

\circ \circ \circ \circ \circ

$$\mathbb{E}[R(s_0, a_0, s_1) + \gamma R(s_1, \mu(s_1), s_2) + \gamma^2 R(s_2, \mu(s_2), s_3) | s_0 = i, a_0 = a]$$



Let  $\mu'$  be a policy such that,

$$\mu'(s) = \begin{cases} \arg \max_a Q_{\mu}(s, a) & ; \quad s = i \\ \mu(s) & ; \quad s \neq i \end{cases}$$

We can show that,

$$v_{\mu^{(i)}}(s) \geq v_{\mu^{(i')}}(s); \forall i' \\ \text{and some } s \in S, \quad v_{\mu^{(i)}}(s) > v_{\mu^{(i')}}(s)$$

This leads to optimal policy  $\mu^*$



## POLICY ITERATION :

Initialization :

Start with a policy  $\mu_0$ .

For  $t = 0, 1, 2, \dots$

- Policy evaluation: compute  $v_{\mu_t}$

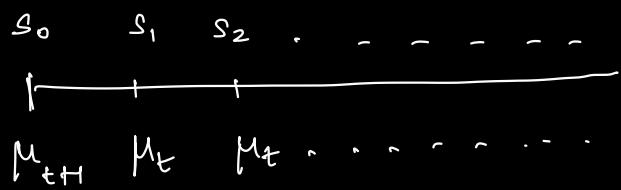
- Policy improvement:

- Compute  $Q_{\mu_t}(s, a) \quad \forall s \in S, a \in A$

-  $\mu_{t+1}^{(i)} = \arg \max_a Q_{\mu_t}(s, a) \quad \forall s \in S$

Repeat the above procedure till,  $\mu_{t+1}^{(i)} = \mu_t^{(i)} \quad \forall i \in S$

At this point Bellman equation gets satisfied.



Claim :

- Following  $\mu'$  for state  $s_0$  at time  $t=0$  and then following  $\mu$  for next time steps  $\rightarrow$  inferior to following  $\mu$  at  $t \geq 2$ .
- $\mu'$  at  $t=0, 1$  and then following  $\mu$  at  $t \geq 2$ .

$$\text{i.e. } Q_{\mu}(i, \mu'(i)) = R(i, \mu'(i)) + \gamma \sum_{j=1}^n P_{ij}(\mu'(j)) V_{\mu}(j) - \textcircled{1}$$

$$= R(i, \mu'(i)) + \gamma \sum_{j=1}^n P_{ij}(\mu'(j)) Q_{\mu}(j, \mu(j))$$

$$Q_{\mu}(i, \mu'(i)) \leq R(i, \mu'(i)) + \gamma \sum_{j=1}^n P_{ij}(\mu'(j)) Q_{\mu}(j, \mu(j))$$

Need to show this !

Following  $\mu'$  twice and then  
following  $\mu$  at  $t \geq 2$

• Expanding (•) :

$$(•) = R(i, \mu'(i)) + \gamma \sum_{j=1}^n P_{ij}(\mu'(j)) \left[ R(j, \mu'(j)) + \gamma \sum_{k=1}^n P_{jk}(\mu'(k)) V_{\mu}(k) \right] - \textcircled{2}$$

$\Rightarrow \textcircled{1} - \textcircled{2}$

$$= V_{\mu}(j) - \left[ R(j, \mu'(j)) + \gamma \sum_{k=1}^n P_{jk}(\mu'(k)) V_{\mu}(k) \right]$$

$$= R(j, \mu(j)) + \gamma \sum_{k=1}^n p_{jk}(\mu(j)) v_{\mu(k)} - R(j, \mu'(j)) \\ - \gamma \sum_{k=1}^n p_{jk}(\mu'(j)) v_{\mu(k)}$$

$$= R(j, \mu(j)) + \gamma \sum_{k=1}^n p_{jk}(\mu(j)) v_{\mu(k)} - \left[ \max_{a \in A} \left( R(j, a) + \gamma \sum_{k=1}^n p_{jk}^{(a)} v_{\mu(k)} \right) \right] \\ \underbrace{\qquad\qquad\qquad}_{\alpha(\mu(j))} \qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{\alpha(a)}$$

Since  $\mu(j) \in A \Rightarrow \alpha(\mu(j)) \leq \max_{a \in A} \alpha(a)$

$$\Rightarrow \textcircled{1} - \textcircled{2} \leq 0$$

$$\Rightarrow \textcircled{1} \leq \textcircled{2}$$

$\Rightarrow$  Taking  $\mu'$  at  $t=0$ , and then  $\mu$  at  $t \geq 2$   
is better than taking  $\mu'$  at  $t=0$  and then  $\mu$  at  $t \geq 1$ .

14/1:

RECAP:

• Bellman Equation:

$$V^*(i) = \max_a \left[ R(i, a) + \sum_{j=1}^n P_{ij}(a) V^*(j) \right]$$

$$\mu^*(i) = \operatorname{argmax}_a \left[ R(i, a) + \sum_{j=1}^n P_{ij}(a) V^*(j) \right]$$

•  $T(V^*) = V^*$

$\underbrace{T : \mathbb{R}^n \rightarrow \mathbb{R}^n}_{\text{contraction operator}}$

• Value Iteration:

$$V_1 = T V_0$$

:

$$\lim_{n \rightarrow \infty} V_n = V^*$$

• Policy Iteration:

1) Policy evaluation:

$$V_\mu = (I - \gamma P_\mu)^{-1} R_\mu$$

2) Policy Improvement:

$$Q_\mu(i, a) = R(i, a) + \gamma \sum_{j=1}^n P_{ij}(a) V_\mu(j)$$

$$\mu^*(i) = \arg \max_a Q_{\mu}(i, a)$$

Let  $T_\mu$  be a Bellman operator for a policy  $\mu$ .

$$T_\mu: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$\therefore T_\mu(v(i)) = R(i, \mu(i)) + \gamma \sum_{j=1}^n P_{i,j}(\mu(i)) v(j)$$

Also,  $T_\mu(v_\mu) = v_\mu$  and  $v_\mu$  is unique.

Value iteration for a single policy.

$$\Rightarrow \lim_{K \rightarrow \infty} T_\mu^K(v_0) = v_\mu.$$

NOTE :  $v^*$  is the fixed point of the entire MDP.

$T$  is for the entire MDP and its fixed point is  $v^*$ .  
 $T_\mu$  is for a fixed policy and its fixed point is  $v_\mu$ .

$$\text{Remember, } T_\mu(v) = R_\mu + \gamma P_\mu \cdot v$$

## Fixed Point of a function :

Consider,  $f : \underbrace{X}_{\text{vector space}} \rightarrow X$

$x^*$  is a fixed point if,

$$f(x^*) = x^*$$

Suppose 'f' is a contraction :

$$\Rightarrow |f(x) - f(y)| \leq \alpha |x-y| ; \alpha \in [0, 1)$$

If 'f' is a contraction then we will have a

unique fixed point.

Proof : Suppose not  $\Rightarrow \exists$  multiple fixed points.

WKT, since 'f' is a contraction,

$$|f(x^*) - f(y^*)| \leq \alpha |x^* - y^*|$$

↙  
fixed points

$$\Rightarrow |x^* - y^*| \leq \alpha |x^* - y^*|$$

- As  $\alpha \in [0, 1] \Rightarrow$  contradiction.

$\therefore$  If ' $f$ ' is a contraction mapping then there exists a unique point.

Ex:

Let  $f(x) = \frac{1}{2} \left( x + \frac{\alpha}{x} \right)$

convex comb.  
of  $x$  &  $\alpha/x$

Neuton's method  
to find the  
square root.

Consider,  $x = \sqrt{\alpha}$  } unique

$$\Rightarrow f(\sqrt{\alpha}) = \frac{1}{2} \left( \sqrt{\alpha} + \frac{\alpha}{\sqrt{\alpha}} \right)$$

$$\underline{\underline{f(\sqrt{\alpha}) = \sqrt{\alpha}}}$$

$$\Rightarrow f^n(x) \rightarrow \sqrt{\alpha}$$



Let  $\alpha = 16$ ,  $x_0 = 2$

$$\Rightarrow \alpha/x_0 = 8 \quad \Rightarrow \quad x_1 = \frac{1}{2}(2 + 8) = 5$$

$$\Rightarrow x_2 = \frac{1}{2}(5 + \frac{16}{5}) = 4.1$$

Claim :

The operator,  $T_\mu : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a contraction operator

w.r.t  $\|\cdot\|_\infty$ -norm.

$$\Rightarrow \|T_\mu(v) - T_\mu(u)\|_\infty \leq \gamma \|v - u\|_\infty$$

Proof :

$$\begin{aligned} |T_\mu(v_1(i)) - T_\mu(v_2(i))| &\leq \left| R(i, \mu) + \gamma \sum_{j \in S} P_{ij}( \mu(i)) v_1(j) \right. \\ &\quad \left. - R(i, \mu) - \gamma \sum_{j \in S} P_{ij}( \mu(i)) v_2(j) \right| \end{aligned}$$

$$\leq \gamma \sum_{j \in S} P_{ij}( \mu(i)) |v_1(j) - v_2(j)|$$

(Triangle's inequality)

$$\text{w.r.t., } |v_1(j) - v_2(j)| \leq \max_{j \in S} |v_1(j) - v_2(j)|$$

$$= \|v_1 - v_2\|_\infty$$

$$\Rightarrow |T_\mu(v_1(i)) - T_\mu(v_2(i))| \leq \gamma \sum_{j \in S} P_{ij}( \mu(i)) \|v_1 - v_2\|_\infty$$

$$\Rightarrow |T_\mu(v_1(i)) - T_\mu(v_2(i))| \leq \underline{\gamma} \|v_1 - v_2\|_\infty \quad ; \quad \forall i$$

$$\max_i |T_\mu(v_1(i)) - T_\mu(v_2(i))| \leq \gamma \|v_1 - v_2\|_\infty$$

$$\Rightarrow \|T_\mu(v_1) - T_\mu(v_2)\|_\infty \leq \gamma \|v_1 - v_2\|_\infty$$

$\therefore T_\mu$  is a contraction operator  $\gamma < 1$ .  
 w.r.t  $\ell_\infty$ -norm and with  
 contraction

Claim : The Bellman operator ' $T$ ' is a contraction operator with respect to  $\ell_\infty$ -norm.  
 with contraction operator  $\gamma < 1$  and w.r.t  $\ell_\infty$ -norm.  
 $\therefore \|T(v_1) - T(v_2)\|_\infty \leq \gamma \|v_1 - v_2\|_\infty ; \gamma \in [0, 1)$

16/2 :

Recap :

Policy evaluation :

Complexity { - Either using direct method i.e  $v_\mu = (\mathbb{I} - \gamma P_\mu)^{-1} R_\mu$   
 is  $O(n^3)$  (OR)

Complexity { - Recursively applying ' $T_\mu$ ' operator with  $v_0 = 0$   
 is  $O(n^2 K)$   $\lim_{K \rightarrow \infty} v_K = v_\mu$  where,  $v_K = T_\mu^K(v_0)$

Claim :

The Bellman operator ' $T$ ' is a contraction operator with contraction operator ' $\gamma < 1$ ' and w.r.t  $\|\cdot\|_\infty$ -norm.

$$\therefore \exists \epsilon \quad \|T(v_1) - T(v_2)\|_\infty \leq \gamma \|v_1 - v_2\|_\infty ; \quad \gamma \in [0, 1)$$

Consider state  $s \in \mathcal{S}$ ,

$$\text{WKT, } T(v_1(s)) = \max_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a) V_1(s') \right]$$

$$T(v_2(s)) = \max_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a) V_2(s') \right]$$

$$\text{Let } a_1^* = \operatorname{argmax}_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a) V_1(s') \right]$$

$$a_2^* = \operatorname{argmax}_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a) V_2(s') \right]$$

$$\Rightarrow T(v_1(s)) - T(v_2(s)) = R(s, a_1^*) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a_1^*) V_1(s') \\ - \left[ R(s, a_2^*) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a_2^*) V_2(s') \right]$$

$$\text{WKT, } R(s, a_2^*) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a_2^*) V_2(s') \geq R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,s'}(a) V_2(s') ; \quad \forall a \in \mathcal{A}$$

Action space

$$\begin{aligned}
\Rightarrow \tau(v_1(s)) - \tau(v_2(s)) &\leq R(s/a_1^*) + \gamma \sum_{s' \in S} P_{s,s'}(a_1^*) v_1(s') \\
&\quad - R(s/a_1^*) + \gamma \sum_{s' \in S} P_{s,s'}(a_1^*) v_2(s') \quad (\because a_1^* \in A) \\
&= \gamma \sum_{s' \in S} P_{s,s'}(a_1^*) (v_1(s') - v_2(s'))
\end{aligned}$$

• WKT  $\forall s' \in S$ ,

$$v_1(s') - v_2(s') \leq \max_{s' \in S} |v_1(s') - v_2(s')| = \|v_1 - v_2\|_\infty$$

$$\begin{aligned}
\Rightarrow \tau(v_1(s)) - \tau(v_2(s)) &\leq \underbrace{\gamma \sum_{s' \in S} P_{s,s'}(a_1^*)}_{= 1} \|v_1 - v_2\|_\infty \\
\tau(v_1(s)) - \tau(v_2(s)) &\leq \gamma \|v_1 - v_2\|_\infty \quad ; \quad \forall s \in S \quad - \textcircled{1}
\end{aligned}$$

• similarly,  $\tau(v_2(s)) - \tau(v_1(s))$

$$\begin{aligned}
\tau(v_2(s)) - \tau(v_1(s)) &\leq R(s/a_2^*) + \gamma \sum_{s' \in S} P_{s,s'}(a_2^*) v_2(s') \\
&\quad - R(s/a_2^*) + \gamma \sum_{s' \in S} P_{s,s'}(a_2^*) v_1(s') \\
&= \gamma \sum_{s' \in S} P_{s,s'}(a_2^*) (v_2(s') - v_1(s')) \\
&\leq \max_{s' \in S} |v_1(s') - v_2(s')| = \|v_1 - v_2\|_\infty
\end{aligned}$$

$$\begin{aligned}
\Rightarrow \tau(v_2(s)) - \tau(v_1(s)) &\leq \underbrace{\gamma \sum_{s' \in S} P_{s,s'}(a_2^*)}_{= 1} \|v_1 - v_2\|_\infty \\
\tau(v_2(s)) - \tau(v_1(s)) &\leq \gamma \|v_1 - v_2\|_\infty \quad ; \quad \forall s \in S \quad - \textcircled{2}
\end{aligned}$$

$$\begin{aligned}
x &\leq y \\
-x &\leq y \\
\Rightarrow |x| &\leq y
\end{aligned}$$

$$\Rightarrow |\tau(v_2(s)) - \tau(v_1(s))| \leq \gamma \|v_1 - v_2\|_\infty ; \quad \forall s \in S \text{ (From ① and ②)}$$

$$\Rightarrow \boxed{\|\tau(v_1) - \tau(v_2)\|_\infty \leq \gamma \|v_1 - v_2\|_\infty}$$

• Comparison of value iteration and policy iteration.

- w.r.t for value iteration,

$$T[V_n(i)] = \max_a \left[ R(i, a) + \sum_{j=1}^n P_{i,j}(a) V_{n-1}(j) \right]; \forall i$$

complexity,  $O(|S|^2 \cdot |A|)$

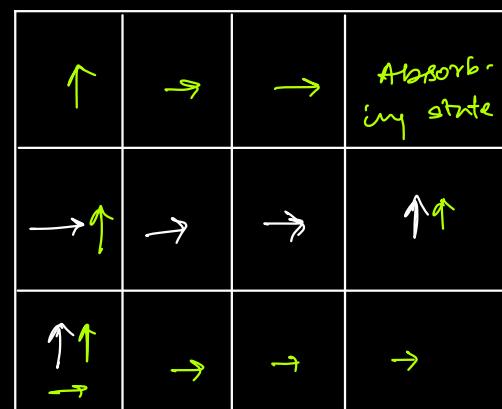
-  $\Rightarrow$  for 'k' steps :  $O(|S|^2 \cdot |A| \cdot K)$

For policy iteration :

- Each step has complexity of  $O(n^3)$
- For 'k' steps it will be  $O(n^3 \cdot k)$

NOTE : Optimal value function is unique where as optimal policy need not be unique.

|   |   |   |   |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |



|            |            |            |            |
|------------|------------|------------|------------|
| $\gamma^3$ | $\gamma^2$ | $\gamma$   | 1          |
| $\gamma^4$ | $\gamma^3$ | $\gamma^2$ | $\gamma$   |
| $\gamma^5$ | $\gamma^4$ | $\gamma^3$ | $\gamma^2$ |

} Unique value function.

### MACHINE REPLACEMENT PROBLEM :

- MDP :  $S := \{1, 2, \dots, n\}$  } set of operations of machine
- Cost function, ' $g^i$ ' is a non-decreasing function of state.
- PTM :
  - i)  $P_{ij} = 0 \quad ; \quad i > j$
  - ii)  $P_{ij} \leq P_{i+1,j} \quad ; \quad i < j$  } Stochastic dominance  
(check!) Yes!

### NOTE :

- A MC with PTM ' $P$ ' is stochastically monotone if,
- $$P_i \leq_s P_j \quad ; \quad \forall i > j$$

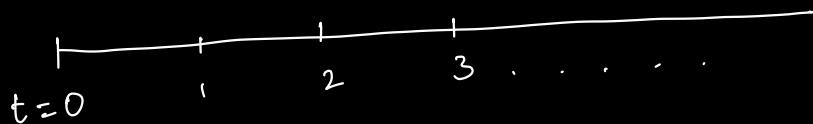
$\underbrace{\quad}_{P}$
- $\Rightarrow$  PMF of  $i^{th}$  row lies below  $j^{th}$  row.

- For a weakly increasing function,  $V: S \rightarrow \mathbb{R}$  with  $s_1 > s_2$  ( $s_1, s_2 \in S$ ) with PTM, P. Then,

$$E[V(s_{t+1}) | s_t = s_1] \geq E[V(s_{t+1}) | s_t = s_2] \text{ iff 'P' is monotone.}$$

- Action: i) Replacement (OR) ii) Continue

- Machine Operation :



- If you replace the machine state 1 for 1 time period. then it will run in

COST FUNCTION :

$$c(i, c, j) = g(i)$$

$$c(i, R, 1) = R + g(i)$$

## PROBABILITY TRANSITION MATRIX :

- $P(1, c, j) = P_{ij}$

- $P(1, R, 1) = 1$

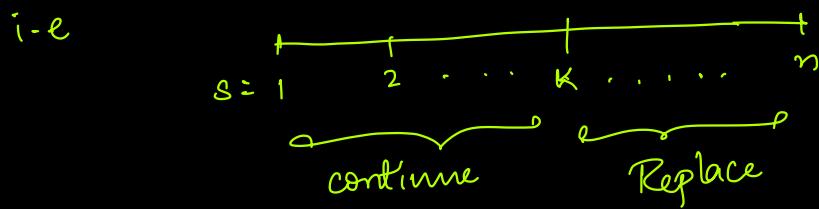
$$P(1, R, \cdot) = 0 ; \forall \cdot \neq 1$$

- $P(\cdot, R, 1) = 1 ;$

## BELLMAN EQUATION :

- $V^*(i) = \min \left\{ q(i) + \gamma \sum_{j=1}^n P_{ij} V^*(j), R + q(1) + \gamma V^*(1) \right\}$

- Intuitively optimal policy is a threshold type policy.



- $V^*(1) < V^*(2) < V^*(3) < V^*(4) \dots \dots \begin{pmatrix} \text{we can show} \\ \text{this is true!} \end{pmatrix}$

Claim:

- 1) Show that  $V^*$  is monotonic (non-decreasing)
- 2) Show that the policy is threshold.

We can prove claim ② assuming claim ① as follows:

$$V^*(i) = \min \left\{ \underbrace{R + g(i)}_{(R)}, \underbrace{g(i) + \sum_{j=1}^n p_{ij} V^*(j)}_{(C)} \right\}$$

since (R) remains same for all  $i > 1$ , if we show w.r.t  $i$  then we are done.

(C) is increasing

w.k.t  $g(i)$  is non-decreasing and also,

$$P_{ij} \leq P_{i+1,j}$$

Now to prove claim 2, consider  $V_0 = 0$ .

$$V_1 = T(V_0)$$

w.k.t,

$$V_1(i) = \min \left\{ R + g(i), g(i) + \sum_{j=1}^n p_{ij} V_1(j) \right\}$$

non-decreasing  $\because g(i)$  is non-decreasing

$$V_2(i) = \min \left\{ R + g(i) + \gamma V_1(i), g(i) + \gamma \sum_{j=1}^n p_{ij} V_2(j) \right\}$$

For some  $i > k$

$$V_2(k) = \min \left\{ R + g(i) + \gamma V_1(i), g(k) + \gamma \sum_{j=1}^n p_{kj} V_1(j) \right\}$$

Since  $g(k) > g(i)$

Also,  $\sum_{j=1}^n p_{kj} V_1(j) - \sum_{j=1}^n p_{ij} V_1(j)$

$$= \sum_{j=1}^k p_{kj} V_1(j) - \sum_{j=1}^i p_{ij} V_1(j)$$

$p_{ij} = 0 ; i > j$   
 $p_{ij} \leq p_{i+1,j} ; i < j$

$$= \sum_{j=1}^i (p_{kj} - p_{ij}) V_1(j) + \sum_{j=i+1}^k p_{kj} V_1(j)$$

$$= \sum_{j=1}^i (\underbrace{p_{kj} - p_{ij}}_{\geq 0}) V_1(j) + \sum_{j=i+1}^k \underbrace{p_{kj} V_1(j)}_{\geq 0}$$

$\Rightarrow V(i)$  is non-decreasing function in  $i$ .

18/2:

- Now when discount factor,  $\gamma = 1$ 
  - For this case we need to have a terminal state,  $s^*$  with  $v(s^*) = 0$  else,  $(I - P_{\mu})$  may not be invertible.
- The MDP problems with  $\gamma = 1$  are called stochastic shortest path problem (OR) Total reward maximization problem.

NOTE :

- For episodic tasks we consider  $\gamma = 1$  and special terminal state
- For continuing task,  $\gamma < 1$ .

STATIONARY RANDOMIZED POLICY :

- $\pi : S \rightarrow \mathbb{P}(A)$

- Here,  $v_{\pi}^{(i)} = \sum_{a,j} \pi(a|i) P_{i,j}(a) R(i,a,j)$   
 $+ \gamma \sum_{a=1}^m \pi(a|i) \sum_{j=1}^n P_{i,j}(a) v_{\pi}^{(j)}$

• Vector Equation :

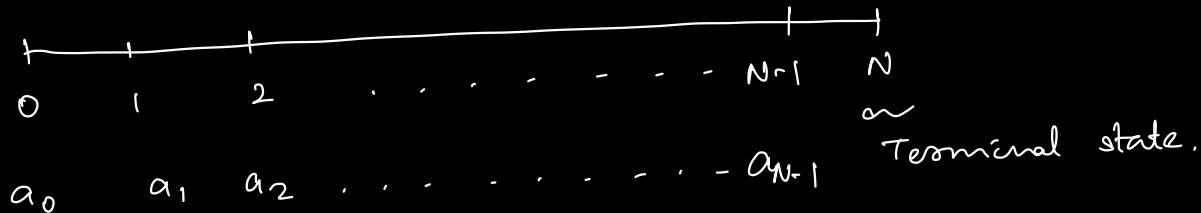
$$V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}$$

where,  $R_{\pi(i)} = \sum_{a \in A} \pi(a|i) \sum_{j=1}^n P_{i,j}(a) R(i,a,j)$

$$P_{\pi(i,j)} = \sum_{a \in A} \pi(a|i) P_{i,j}(a)$$

• For a finite MDP, we will always have an optimal stationary deterministic policy.

### FINITE HORIZON PROBLEM :



• Ex : Chess game :

Policies : 1) Play Aggressively

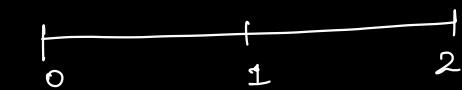
$$\text{Result} = \begin{cases} \text{win w.p } p_w \\ \text{lose w.p } 1-p_w \end{cases}$$

2) Play defensively :

$$\text{Result} = \begin{cases} \text{Draw} & w.p \quad P_d \\ \text{Loose} & w.p \quad 1 - P_d \end{cases}$$

a 3-match series chess game :

Consider



If we are tied after playing game '0' and

then for game '1' we play aggressively.

'1' (OR)  
The player who plays aggressively wins.

State Space :

Let  $S = [w_1, w_2]$

$\downarrow$                            $\curvearrowright$   
 # of times                  # of times  
 we win                  opponent player  
 wins.

Initial state :  $s_0 = (0, 0)$

state at time '1'

$$s_1 \in \{(1, 0), (0, 1), (0.5, 0.5)\}$$

$$\text{state at time '2'} \quad s_2 \in \{(2, 0), (1, 1), (1.5, 0.5), (0, 2), (0.5, 1.5)\}$$

In finite horizon problem, we need not consider PTM to be stationary. Also reward function need not be stationary.

i.e. PTM can vary w.r.t time.

How to solve the finite horizon problem?  
- We use the Backward induction technique.

### BACKWARD INDUCTION TECHNIQUE :

We define,

Is this our assumption?

$$V_2(s) = \begin{cases} 1 &; w_1 > w_2 \\ p_n &; w_1 = w_2 \\ 0 &; w_1 < w_2 \end{cases}$$

*Is this after playing game '2'?*

How we find  $V_1$ ?

$$\text{WKT, } V_1(s_1) = \max_a \left( R(s, a) + \mathbb{E} \left[ J_2(s_2) \right] \right)$$

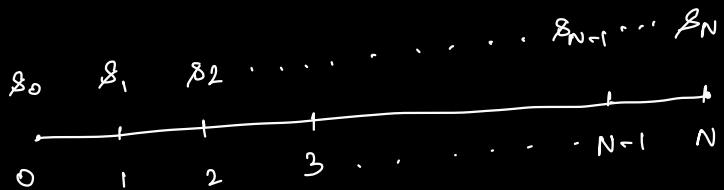
$$u(s_1) = \arg \max_a \left( R(s, a) + \sum_{s_2} P_{s_1, s_2} V_2(s_2) \right)$$

$$\text{Why, } J_0(s_0) = \max_{a_0} \left[ R(s_0, a_0) + \sum_{s_1} P_{s_0, s_1} V_1(s_1) \right]$$

- Here,  $\mu_0 : S \rightarrow A$ ,  $\mu_1 : S \rightarrow A$ ,  $\mu_2 : S \rightarrow A$   
 These policies need not be same.

28/2:

## FINITE HORIZON PROBLEM:



- $J_0(s_0) = \mathbb{E} [R_0 + R_1 + R_2 + \dots + R_{N-1} + R_N]$   
 Expectation is w.r.t actions  $a_0, a_1, \dots, a_{N-1}$
- $J_N(s_N) = R_N$

- $R_t = R_t(s_t, a_t, s_{t+1})$
- $P_t(s_{t+1} | s_t, a_t)$  : Can be time dependant

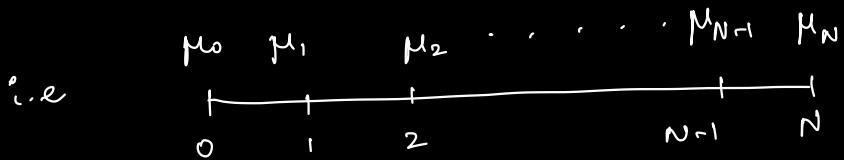
$$\begin{aligned}
 J_{N-1}(s_{N-1}) &= \max_{a_{N-1}} \mathbb{E} \left[ R(s_{N-1}, a_{N-1}, s_N) + J_N(s_N) \right] \\
 &= \max_{a_{N-1}} \left( \sum_{s_N \in S} \left( R(s_{N-1}, a_{N-1}, s_N) + J_N(s_N) \right) P(s_N | s_{N-1}, a_{N-1}) \right)
 \end{aligned}$$

- In general,

$$J_K(s_k) = \max_{a_k} \left( R(s_k, a_k) + \sum_{s_{k+1}} P^s_{s_k, a_k} J_{K+1}(s_{k+1}) \right)$$

- For finite horizon problem the policy is time

dependent.



NOTE :

Even when the reward and probability transitions are time independent this doesn't guarantee that stationary policy is optimal. ( $\because$  the horizon itself is finite)

CHESS PROBLEM :

Let  $x = w_1 - w_2$  (Net score)

$$\Rightarrow J_2(x) = \begin{cases} 1 & ; \text{ if } x > 0 \\ 0 & ; \text{ if } x < 0 \end{cases}$$

When the final score is

(1,1) we enter sudden-death i.e players

keep playing until

someone wins.

$\xrightarrow{\text{Why makes sense}}$  when both the players have won equal games?

$\rightarrow$  playing boldly will end the sudden-death match at once.

$$\Rightarrow J_2 \left( \underbrace{(2, 0)}_{(w_1, w_2)} \right) = 1$$

Return due to  
playing aggressively

$$\therefore J_1 \left( (1, 0) \right) = \max \left\{ \begin{array}{l} p_w J_w(2, 0) + (1-p_w) J_w(1, 1), \\ p_d J_2(1.5, 0.5) + (1-p_d) J_2(1, 1) \end{array} \right\}$$

playing defensively

$$= \max \left\{ \underbrace{p_w + (1-p_w) p_w}_{a}, \underbrace{p_d + (1-p_d) p_w}_{b} \right\}$$

$$\text{WKT, } p_d > p_w$$

$$\Leftrightarrow p_w - p_w^2 < p_d + p_w - p_d p_w$$

$$\Rightarrow p_w - p_w^2 - p_d + p_d p_w < 0$$

$$\Rightarrow p_w - p_d + p_w(p_d - p_w) < 0$$

$$\Rightarrow \underbrace{(p_w - p_d)}_{\leq 0} \underbrace{(1 - p_w)}_{> 0} < 0$$

$$\Rightarrow a \leq b$$

⇒ For  $s_1 = (1, 0)$  play defensively.

In general,

$$J_t(x_t) = \max_{w_1, w_2} \left\{ \begin{array}{l} p_d J_{t+1}(x_t) + (1-p_d) J_{t+1}(x_{t-1}) \\ p_w J_{t+1}(x_{t+1}) + (1-p_w) J_{t+1}(x_{t-1}) \end{array} \right\}$$

$$\mu_t(x_t) = \begin{cases} \text{play bold ; if } p_d J_{t+1}(x_t) + (1-p_d) J_{t+1}(x_{t-1}) \\ & < p_w J_{t+1}(x_{t+1}) + (1-p_w) J_{t+1}(x_{t-1}) \\ \text{play defensively ; otherwise} \end{cases}$$

$$\Rightarrow J_1(0-1) = \max \left\{ \begin{array}{l} p_d J_2(0-5-1-5) + (1-p_d) J_2(0-2) \\ p_w J_2(1-1) + (1-p_w) J_2(0-2) \end{array} \right\}$$

$$= \max \left\{ p_d(0) + (1-p_d)(0), p_w \cdot p_w + (1-p_w)(0) \right\}$$

$$= \max \{ 0, p_w^2 \}$$

$$\Rightarrow J_1(0-1) = p_w^2 ; \quad \mu_1(0-1) : \text{Play aggressively}$$

$$\text{If } J_1(0.5-0.5) = \max \left\{ \begin{array}{l} p_d J_2(1-1) + (1-p_d) J_2(0.5-1.5), \\ p_w J_2(1.5-0.5) + (1-p_w) J_2(0.5-1.5) \end{array} \right\}$$

$$= \max \left\{ p_d \cdot p_w, p_w \right\}$$

$$= p_w \quad (\because p_d < 1)$$

$\Rightarrow \mu_1(0.5-0.5) : \text{Play aggressively}$

$$\begin{aligned} \therefore J_0(0-0) &= \max \left\{ \underbrace{p_d J_1(0.5-0.5)}_{a} + \underbrace{(1-p_d) J_1(0-1)}_{b}, \right. \\ &\quad \left. \underbrace{p_w J_1(1-0) + (1-p_w) J_1(0-1)}_{b} \right\} \\ &= \max \left\{ p_d \cdot p_w + (1-p_d) p_w^2, \right. \\ &\quad \left. p_w \left( p_d + (1-p_d)p_w \right) + (1-p_w) p_w^2 \right\} \end{aligned}$$

$$\Rightarrow \cancel{p_d p_w} + \cancel{(1-p_d) p_w^2} \quad ( ) \quad \cancel{p_w p_d} + \cancel{p_w^2 - p_w^2 p_d} + \cancel{p_w^2 - p_w^3}$$

$$0 \quad (<) \quad \underbrace{p_w^2}_{>0} \underbrace{(1-p_w)}_{>0}$$

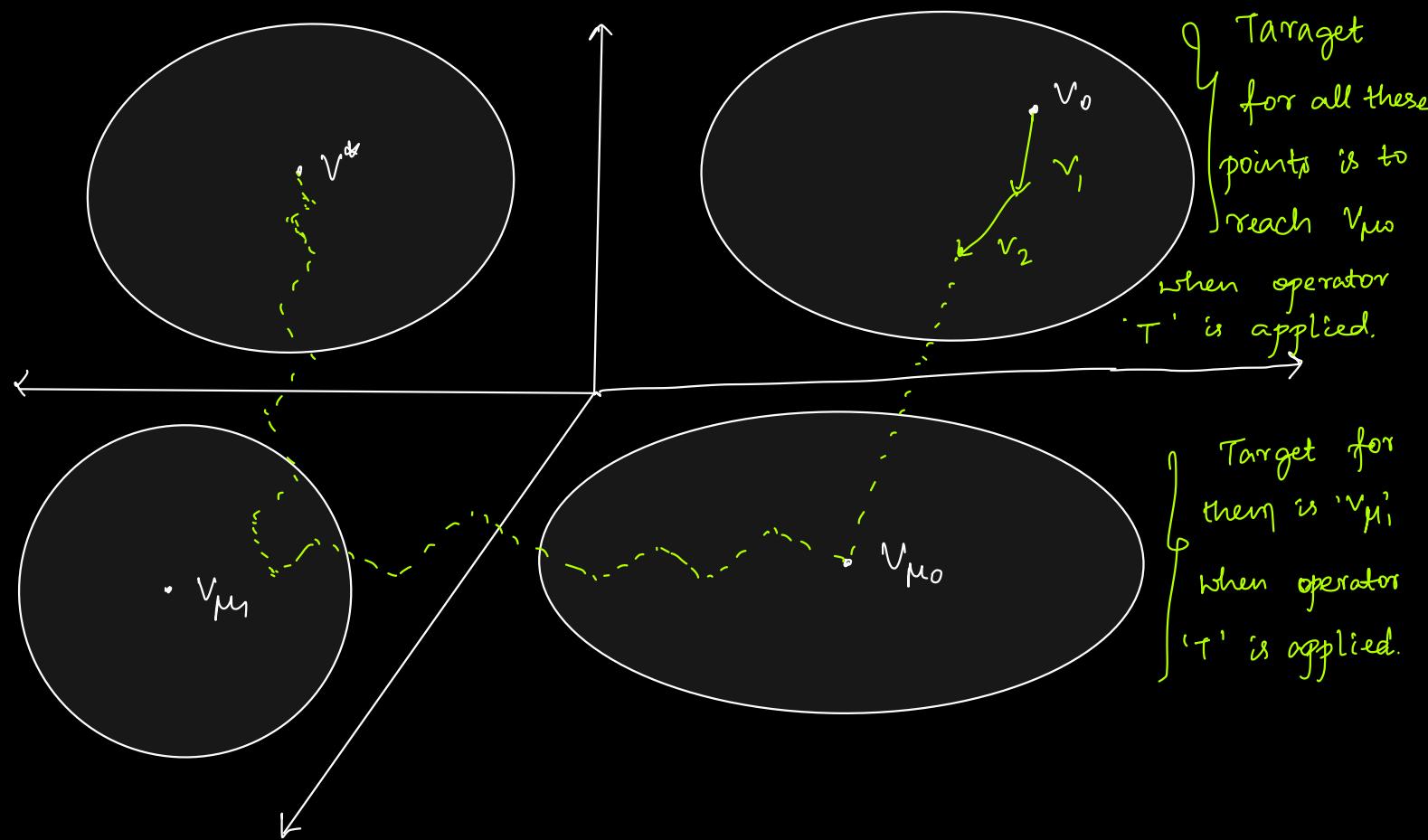
$\Rightarrow \mu_0(0-0) : \text{Play aggressively.}$

Conclusion :

- Play aggressively in the starting match
- Play aggressively if we are ahead in the next matches.

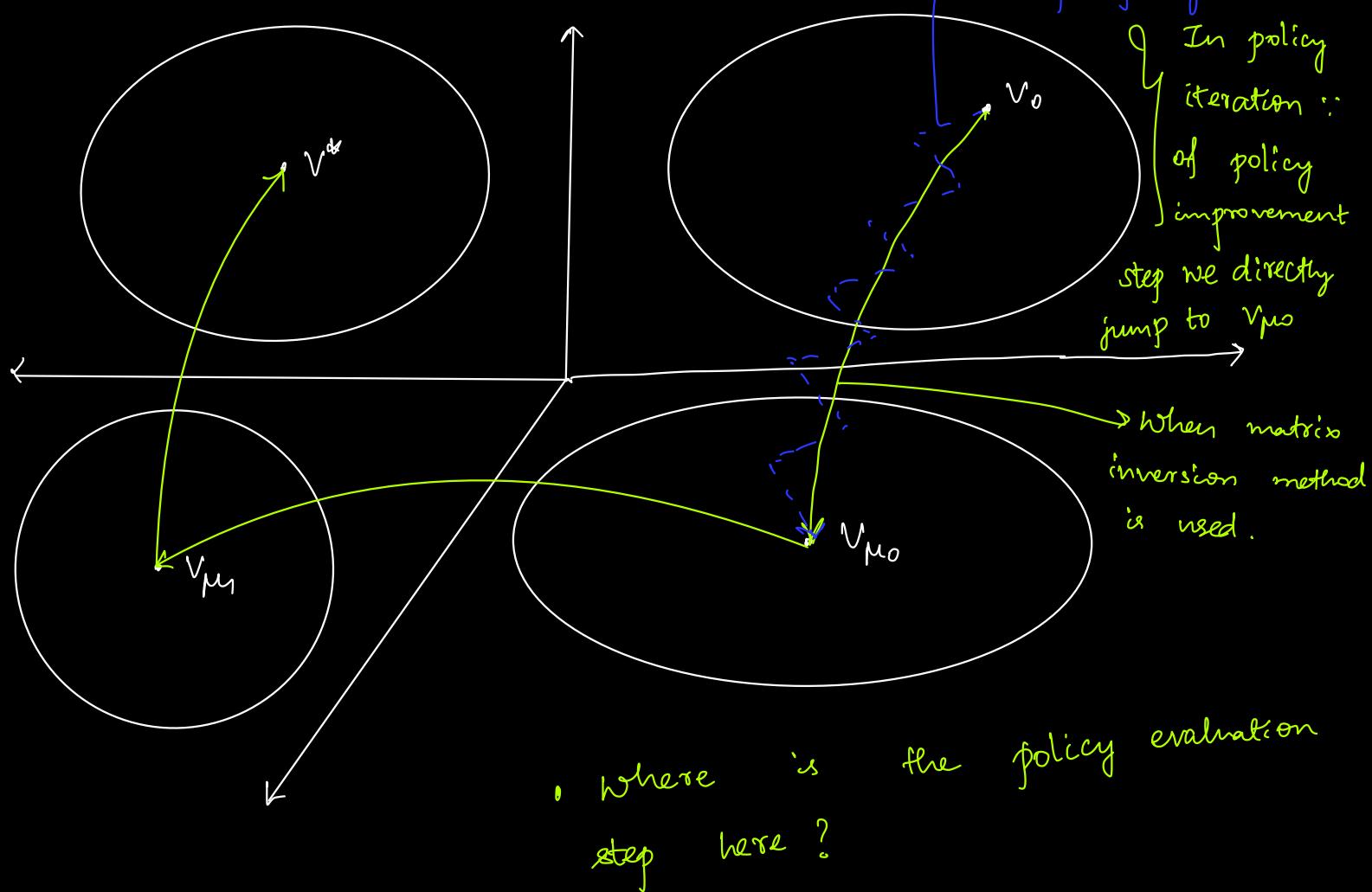
### PICTORIAL REPRESENTATION OF VALUE AND POLICY ITERATION:

• Value Iteration :  $T^n$



## • Policy Iteration :

$\mathbb{R}^n$



• Where is the policy evaluation step here?

• Why the regions are contiguous?

{ why up  
both these  
regions have  
same  $\mu$  as  
greedy policy.

$$R(s, \mu(s)) + \sum_{s'} P_{s,s'} \mu(s) J(s)$$

$$> R(s, a) + \sum_{s'} P_{s,s'} a J(s); \quad \forall a \neq \mu(s)$$

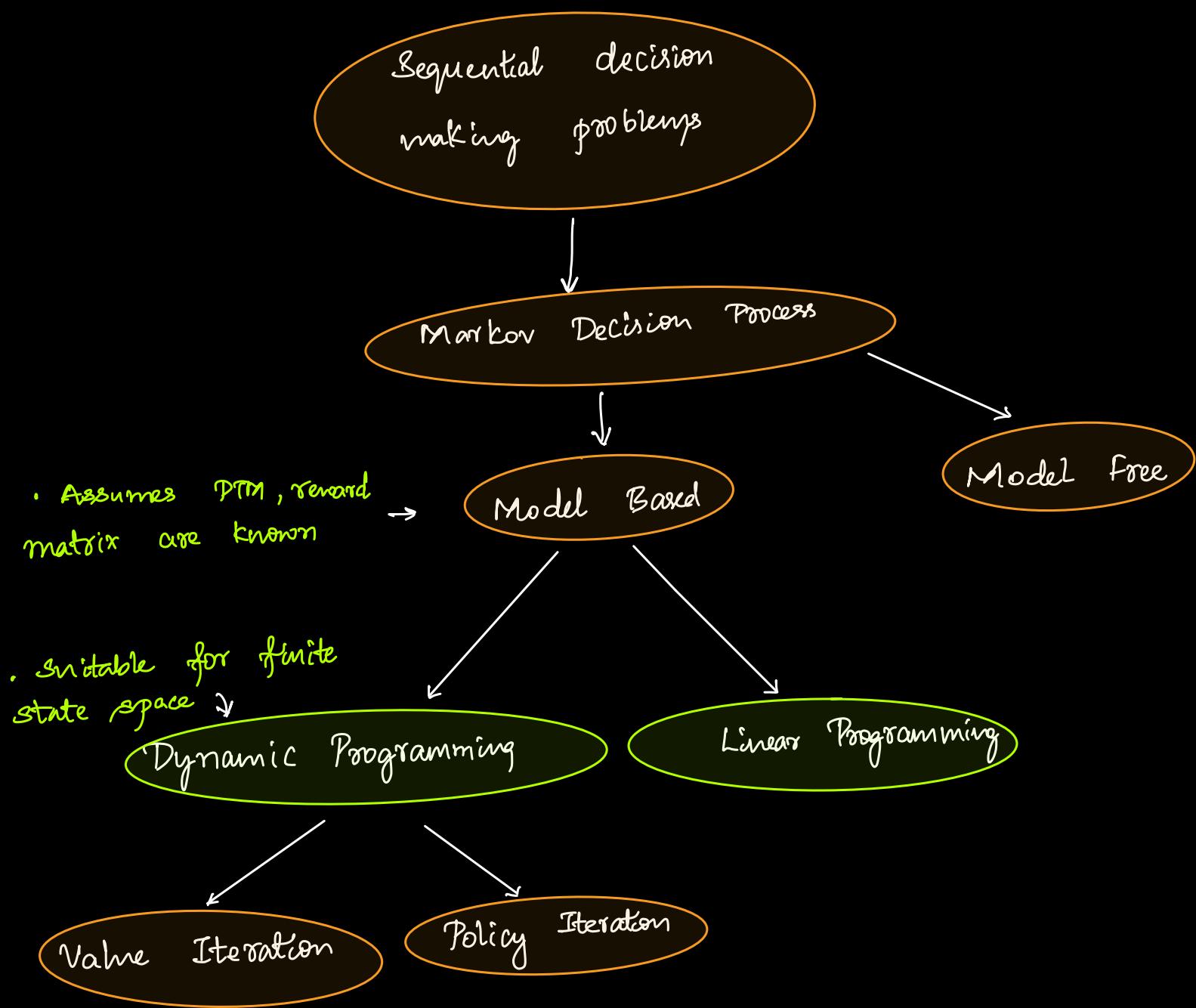
} Also true  $\forall s \in S$

For all these  $\Rightarrow$  There are  $|A|-1$  constraints.  
 $J(s), \mu(s)$  will be greedy (which are linear constraints)

So the region will be a polytope with the region of interest being the intersection of the region formed due to linear constraints.

2/3:

BIG PICTURE OF RL :



## Model-free methods

• Reward and PTM  
are unknown

- Suitable for smaller state spaces

Without function approximation / Tabular methods

With function approximation

- Model Free Methods can be classified into :

### i) Monte Carlo Methods

depends on a particular trajectory (path) to reach the end state

### ii) Temporal difference methods

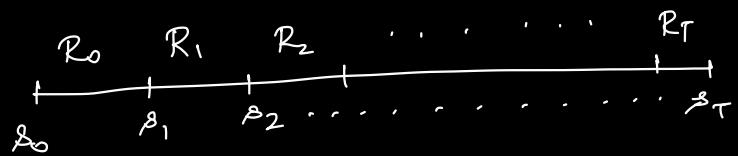
- Why model free methods instead of estimating PTM and reward matrix.
  - We can do this but this might be a tedious task.

- The main steps in RL based methods is
  - Prediction (similar to policy evaluation)
  - Control (similar to policy improvement)
- For prediction we use either MC or TD method.

**NOTE :**

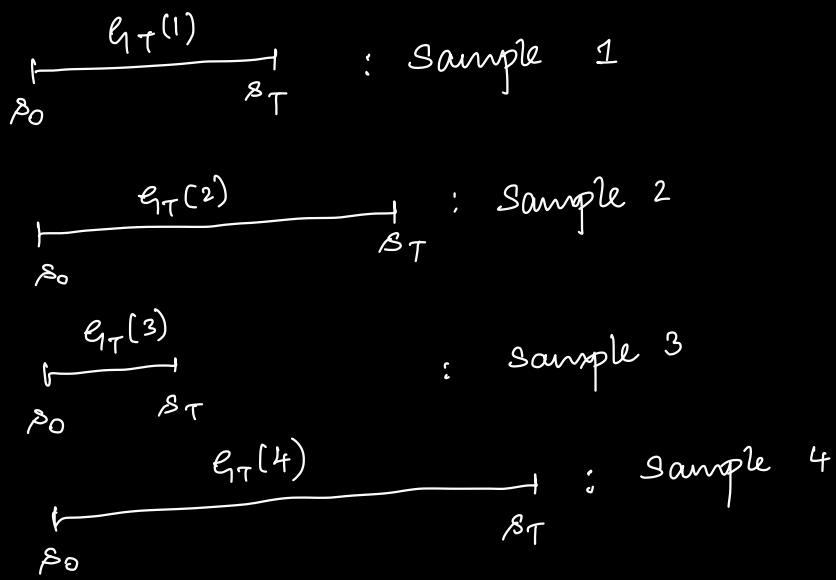
- In all the settings that we are going to deal from now on, we consider shortest stochastic path objective criterion (i.e total reward criterion)
- $\Rightarrow$  Assumption of terminal state.

wk<sup>T</sup>, for a particular trajectory :



$$\Rightarrow v_{\mu^{(s_0)}} = \mathbb{E} \left[ \sum_{t=0}^T R_t \right]$$

Can't be computed since  
the dynamics of the model  
is unknown.



$$\Rightarrow V_M(s_0) = \sum_{i=1}^4 \frac{e_{1:T}(i)}{4}$$

In general for any state,  $s$ ,

$$V_M(s) = \sum_{i=1}^N \frac{e_{1:T}(i) \mid s_t = s : s_T]}{N}$$

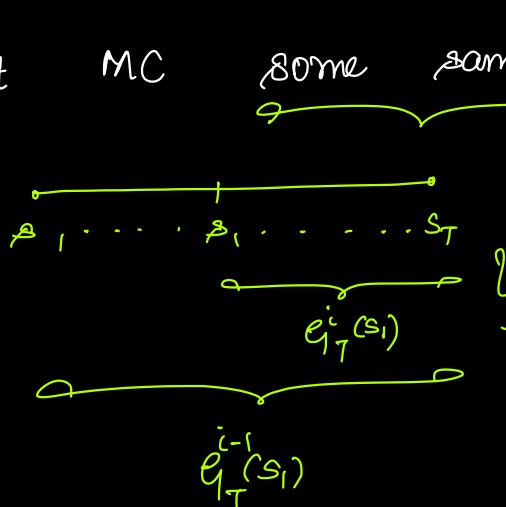
Return obtained  
from state ' $s$ '  
till the terminal state ' $s_T$ '.

It is possible that a state ' $s$ ' can be visited more than once in single trajectory itself. Based on the usage of return from this intermediate visit to the terminal we have 2 methods:

1) First visit Monte-Carlo

2) Every visit Monte-Carlo

$\curvearrowright$  Samples are correlated

- In First visit MC each sample is i.i.d.  
 a trajectory starting with our state of interest till the end of the episode.
  - In Every visit MC some samples are correlated.  


This return is dependent on the trajectory of the return  $e_{1:T}^{(i-1)}(s_1)$
  - We can see that, first visit MC is unbiased but has high variance. ( $\because$  Total variance  $\propto \frac{1}{n} \rightarrow$  # of samples
  - Also, every visit MC is biased and the variance is low when compared with first visit MC.
- By how much the variance is low depends on how many times we visit the same state in a trajectory  
 $(\because$  Total variance  $\propto \frac{1}{m})$  where  $m > n$

- How do we improve the policy?
  - We need  $Q_{\mu}(s, a)$  for policy improvement.
  - Since we only have policy ' $\mu$ ' known to us, we may not be able to estimate  $Q_{\mu}(s, a)$   $\forall a \neq \mu(s)$ .

• Dark Analogy for Bias-Variance Trade-off:

$E[\hat{\theta}]$

$\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$

$\hat{\theta}$

$\theta = E[\hat{\theta}]$

Bias is non-zero  
but variance is less

Bias is zero  
but high variance

- Hence we use some exploration strategy.

i.e.

$$\mu = \begin{cases} \mu & \text{w.p } 1-\epsilon \\ \text{other actions} & \text{w.p } \epsilon \end{cases}$$

} on-policy  
MC control  
(Randomized Policy)

- $\mu^*(s) = \arg \max_a Q_{\mu}(s, a)$

4/3 :

Recap:

- MC control is a full RL method. It consists of:
  - ▷ Policy evaluation - stationary randomized policy ( $\epsilon$ -greedy)

On-policy control

{ Evaluate  $\hat{Q}_\pi$  using

- 1) First visit MC
- 2) Every visit MC

ii) Policy improvement

$$\mu'(s) = \underset{a}{\operatorname{argmax}} Q_\mu(s, a)$$

Ex:

Game of Black Jack:

. Dealer

[up] [down] [down]

Agent :

Goal to have score close to 21

Action : Hit / Stick

• State space : (3 dimensional)

• Sum (11 - 20) of the player

• Face up card of dealer :  $\{1, \dots, 10\}$

• Usable ace or not :  $\{0, 1\}$

$$\Rightarrow |S| = \underline{\overline{200}}$$

• Action = {Hit, Stick}

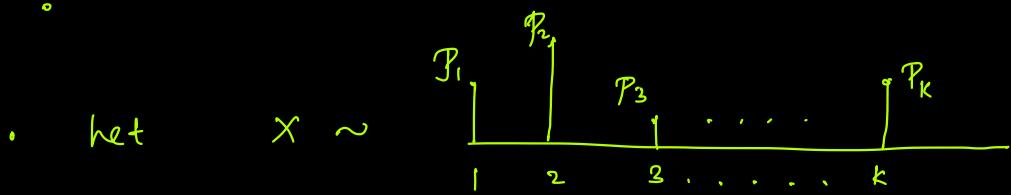
OFF POLICY MONTE CARLO :

• In off-policy MC our goal is to learn the value function corresponding to target policy  $\pi^*$ .

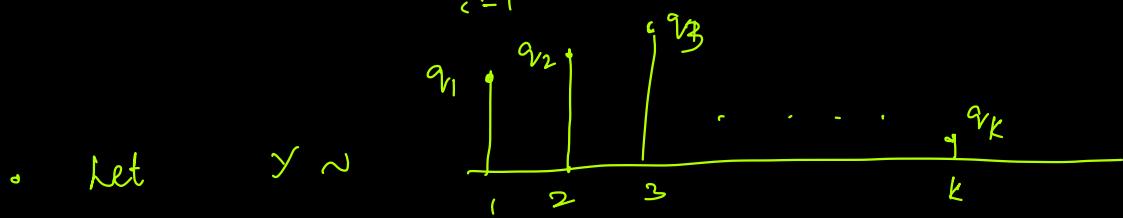
However the actions corresponding to only behavior policy is allowed.

• How to estimate  $V_{\pi^*}$  by following policy  $\pi$ ?  
- Use importance sampling.

NOTE :



$$\mathbb{E}X = \sum_{i=1}^k i \cdot p_i \quad \text{by our interest}$$



$$\mathbb{E}Y = \sum_{i=1}^k i \cdot q_i$$

- We want to estimate  $\mathbb{E}X$  whose  $p_i$ 's are known. However we are restricted to using only samples of  $Y$  that follows some other distribution.

- For each realization of  $Y$  we weight it by

$$(p_i/q_i)$$

$$\Rightarrow = \sum_{i=1}^n i \cdot \frac{p_i}{q_i} \cdot q_i$$

$$\mathbb{E}X = \sum_{i=1}^n i \cdot p_i$$

In general, let  $x \sim f(x)$ ,  $y \sim g(x)$

$$\text{WKT, } E_x = \int_{-\infty}^{\infty} x \cdot f_x(x) dx$$

$$= \int_{-\infty}^{\infty} x \cdot f_x(x) \cdot \frac{g_x(x)}{g_x(x)} \cdot dx$$

$$= \int_{-\infty}^{\infty} x \cdot \frac{f_x(x)}{g_x(x)} \cdot g_x(x) dx$$

$$E_x = \mathbb{E}_{x \sim g_x} \left[ x \cdot \frac{f_x}{g_x} \right]$$

ORDINARY IMPORTANCE

SAMPLING:

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n x_i \cdot \frac{f(x_i)}{g(x_i)} \longrightarrow E_{x \sim f_x}[x] \quad \begin{cases} \text{sampled using dist } 'g'. \\ \text{Importance factor: Specifies how important } 'x' \\ \text{is while finding the expectation w.r.t } 'f_x'. \end{cases}$$

WEIGHTED  
(Normalized)

IMPORTANCE

SAMPLING:

$$\frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i} \longrightarrow E_x \quad \begin{cases} \text{Biased but low variance} \end{cases}$$

$$\text{where, } w_i = \frac{f_i}{g_i}$$

• conditions on  $q_i$  are :

$$i) q_i > 0 \quad \text{And} \quad p_i > 0$$

• when  $q_i$ 's are close to '0' then the weights  $w_i$ 's will be very large resulting in instability in the (oscillations in the estimation)

case of Ordinary IS.

• In weighted IS the denominator  $\sum_{i=1}^m \frac{p_i}{q_i}$  will take care of this. This results in ratios  $\tilde{w}_i$  being very small.

$$\sum_{i=1}^m \frac{1}{\left(\frac{p_i}{q_i}\right)} = \sum_{i=1}^m \frac{p_i}{q_i} x_i$$



9/3:

Recap:

- Monte-Carlo Prediction:

Incremental update of Monte-Carlo:

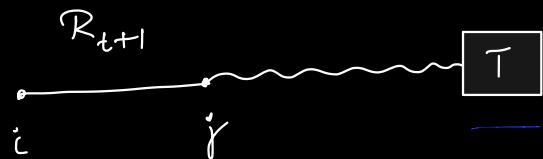
$$V_{t+1}^{(i)} = V_t^{(i)} + \underbrace{\alpha_n (e_t - V_t^{(i)})}_{\frac{1}{n+1}}$$

- WKT,  $V_{\mu}^{(i)}$ : Target value

$V_t^{(i)}$ : Estimate of value function at time  $t$

- New estimate = Current Estimate +  $\alpha$  (sample - current estimate)

TEMPORAL DIFFERENCE PREDICTION:



TD update rule

$$V_{t+1}^{(i)} = V_t^{(i)} + \alpha_n (R_{t+1} + \gamma V_t^{(j)} - V_t^{(i)})$$

- As  $t \rightarrow \infty$ ,  $V_t^{(i)} \rightarrow V_{\mu}^{(i)}$   $\forall i$

• Convergence :

$$v_{t+1}(i) = v_t(i) + \underbrace{\alpha_t \delta_t}_{\text{Temporal difference}}$$

$$\text{for convergence: } \sum_{t=0}^{\infty} \alpha_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

$$\Rightarrow v_t \xrightarrow{t \rightarrow \infty} v_{\mu}$$

Let  $T_{\mu}$  be a contraction operator for a fixed policy  $\mu$  with contraction factor ' $\delta$ '.

$$\Rightarrow v_{t+1} = T_{\mu} v_t \quad \text{and} \quad v_t \xrightarrow{t \rightarrow \infty} v_{\mu}$$

Now consider,

$(1-\alpha) I + \alpha T_{\mu}$

$\overbrace{\hspace{10em}}$  Is this a contraction

$$\text{w.r.t., } \|((1-\alpha)I + \alpha T_{\mu})x - ((1-\alpha)I + \alpha T_{\mu})y\|$$

$$= \| (1-\alpha)I(x-y) + \alpha T_{\mu}(x-y) \|$$

$$\leq \| (1-\alpha)I(x-y) \| + \alpha \| T_{\mu}(x-y) \|$$

$$\leq \underbrace{(1-\alpha) + \alpha \delta}_{= 1 - \alpha(1-\delta)} \|x-y\| < 1$$

$\Rightarrow (1-\alpha)I + \alpha T_\mu$  is a contraction operator with contraction factor  $(1-\alpha)I + \alpha S$ .

Comparing

$$1 - \alpha(1 - \gamma) \geq \gamma$$

$$1 - \alpha + \alpha\gamma - \gamma \geq 0$$

$$1 - \alpha - \gamma(1 - \alpha)\gamma \geq 0$$

$$\underbrace{(1-\alpha)}_{>0} \underbrace{(1-\gamma)}_{>0} \geq 0$$

$\Rightarrow T_\mu$  is a better contraction operator

compared to  $(1-\alpha)I + \alpha T_\mu$ .

$\underbrace{(1-\alpha)I + \alpha T_\mu}_{\text{slow contraction}}$

We can see that for TD prediction:

$$V_{t+1}(i) = (1-\alpha)V_t(i) + \alpha(R_t(i) + \gamma V_t(j))$$

Also,

$$V_{t+1}(i) = \underbrace{(1-\alpha)I + \alpha T_\mu}_{X_{t+1}} V_t(i)$$

$$X_{t+1} = T(x_t)$$

$$= (1-\alpha) V_t(i) + \alpha T_{\mu}(V_t(i))$$

$$V_{t+1}(i) = (1-\alpha)V_t(i) + \alpha \mathbb{E}_{a \sim \mu(i)} \left[ R_t(i, a) + \gamma \sum_{j \in S} p(i, j) V_t(j) \right]$$

- Using just a sample in the place of expectation, we have :

$$V_{t+1}(i) = (1-\alpha)V_t(i) + \alpha \left[ R_t(i, \mu(i), j) + \gamma V_t(j) \right]$$

↑  
map to stochastic version of value iteration  
of fixed policy.

### TEMPORAL DIFFERENCE CONTROL :

- Estimate Q-value function and update based on estimated Q-values.

- Update rule for state-action value is,

$$Q_{t+1}(i, a) = Q_t(i, a) + \alpha (R_{t+1} + \gamma Q_t(j, b) - Q_t(i, a))$$

- w.r.t, a deterministic policy can't be used ( $\because$  all state-action values won't get updated)

• Let ' $\mu$ ' be a randomized policy ( $\varepsilon$ -greedy)

corresponding to deterministic policy  $\pi$ .

} SARSA  
algorithm

• Update,  $Q_\pi(i, a) \forall i, a$

• Update,  $\pi^*(i) = \arg \max_a Q_t(i, a)$

• Comparison of TD learning and MC learning :

• TD learning :

Advantages: 1) Updating is online  
2) can be used directly for continuing task  
3) very helpful if underlying problem has Markov property.  
( $\because$  we use recursion)

• Example :

• Prediction :

Episode 1 : A 0 B 0 6 : B 1

2 : B 1 7 : B 1

3 : B 1 8 : B 1

4 : B 0

5 : B 1

MC

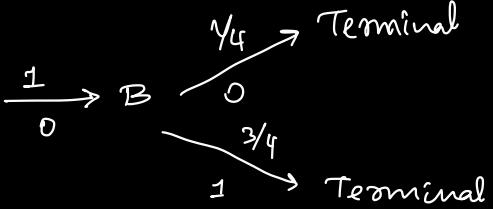
Method :

$$V(A) = 0$$

$$V(B) = 6/8$$

This is  
the model  
the TD method  
has assumed

TD Method :



$$\Rightarrow \hat{V}(B) = 3/4$$

$$\hat{V}(A) = 3/4$$

Why:

Recap:

WKT in value iteration with  $v_0 = 0$ ,

$$- v_k = T_\mu v_{k-1} \Rightarrow v_k \rightarrow v_\mu$$

$$- v_k = \tau v_{k-1} \Rightarrow v_k \rightarrow v^*$$

$$- \text{Also, } ((1-\alpha)I + \alpha T_\mu) v_k \rightarrow v_\mu$$

TD Learning:

$$v_{t+1}(i) = v_t(i) + \alpha_t [R_{t+1} + \gamma v_t(j) - v_t(i)]$$

$$= (1-\alpha_t)v_t(i) + \alpha_t (R_{t+1} + \gamma v_t(j))$$

If multiple samples are allowed then we have,

$$V_{t+1}(i) = (1 - \alpha_t) V_t(i) + \alpha_t \mathbb{E}[R_{t+1} + \gamma V_t(j)]$$

\$\xrightarrow{\text{MVR}}\$ to  $(1 - \alpha_t) I + \alpha_t T_M V_t$

- Here,  $\sum \alpha_t = \infty$ ,  $\sum \alpha_t^2 < \infty$   
 $(\text{Ex: } V_t)$   $\xrightarrow{\quad \downarrow \quad}$   
 To make sure that  
 we can reach anywhere  
 starting from anywhere  
 To make sure that  
 there is no large  
 deviation b/w consecutive  
 time-steps.

- For  $\alpha_t = \gamma_t$

$$\sum \alpha_t = 1 + \gamma_2 + \gamma_3 + \dots = O(\log t)$$

$$\sum \alpha_t^2 = 1 + \gamma_1^2 + \gamma_2^2 + \dots \rightarrow \pi^2/6$$

- $\sum \alpha_t^2 < \infty$  is to cancel the noise.

$\Rightarrow$  For later points we go slowly i.e.  $\alpha_t \downarrow$ .

SARSA is a on-policy algorithm as the update is based on the current policy itself.

Will  $(1-\alpha)I + \alpha T \rightarrow V^*$  ?

wkT,  $V_{t+1} = ((1-\alpha)I + \alpha T)V_t \rightarrow V^* \text{ (Yes!)}$

Consider the stochastic version of the above update:

$$\Rightarrow V_{t+1}(i) = (1-\alpha)V_t(i) + \alpha \left( \max_a \left( R(i, a) + \sum_{j=1}^n P_{ij}(a) V_t(j) \right) \right)$$

$\downarrow$

$V^*$

Here we cannot take a particular sample :: of the max operator.

Now consider,

$$Q^*(i, a) = R(i, a) + \gamma \sum_{j \in S} P_{ij}(a) \max_b Q^*(j, b)$$

$\therefore$  The corresponding stochastic version is,

$Q$ -Bellman equation } 
$$Q_{t+1}(i, a) = (1-\alpha)Q_t(i, a) + \alpha \left[ R(i, a) + \max_b Q_t(j, b) \right]$$

$$Q_{t+1}(i, a) = Q_t(i, a) + \alpha \left[ R(i, a) + \max_b Q_t(j, b) - Q_t(i, a) \right]$$

$$Q_t(i, a) \longrightarrow Q^*(i, a)$$

- For the  $Q$ -learning the behavior policy is any randomized policy (Ex:  $\epsilon$ -greedy policy)
- The target policy here is the stationary deterministic policy.
- $Q$ -learning is a off-policy algorithm.

## SUMMARY :

- | 1) MC control                                                                                                                                                                                                                                                       | 2) SARSA                                                                                                                                                                                    | 3) $Q$ -learning                                                                           |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------|
|                                                                                                                                                                                                                                                                     | For small state-action space<br>when model info. is unknown                                                                                                                                 |                                                                                            |
| <ul style="list-style-type: none"> <li>Uses PI as the dynamic programming principle</li> <li>No bootstrapping</li> <li>Need to wait until the end of episode</li> <li>Mainly for episodic task</li> <li>Reliable even if Markov property is not present.</li> </ul> | <ul style="list-style-type: none"> <li>Uses PI to improve policy but also VI to estimate the policy.</li> </ul>                                                                             | <ul style="list-style-type: none"> <li>Uses VI of max. term in the update step.</li> </ul> |
|                                                                                                                                                                                                                                                                     | ↗<br>↗                                                                                                                                                                                      |                                                                                            |
|                                                                                                                                                                                                                                                                     | <ul style="list-style-type: none"> <li>Uses bootstrapping</li> <li>Online</li> <li>Can be used for continuous and episodic tasks</li> <li>Reliable if Markov property is present</li> </ul> |                                                                                            |

W3:

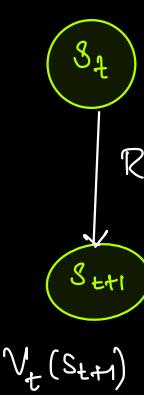
Recap:

• TD prediction :

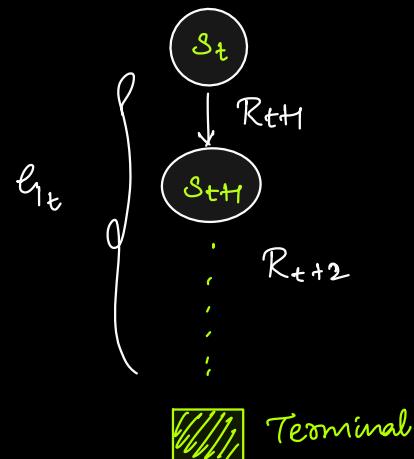
$$V_{t+1}(i) = V_t(i) + \alpha \left[ R_{t+1} + \gamma V_t(j) - V_t(i) \right]$$

$\underbrace{R_{t+1} + \gamma V_t(j) - V_t(i)}$   
 $\delta_t$  (TD error)

• TD Backup:



• MC Backup:



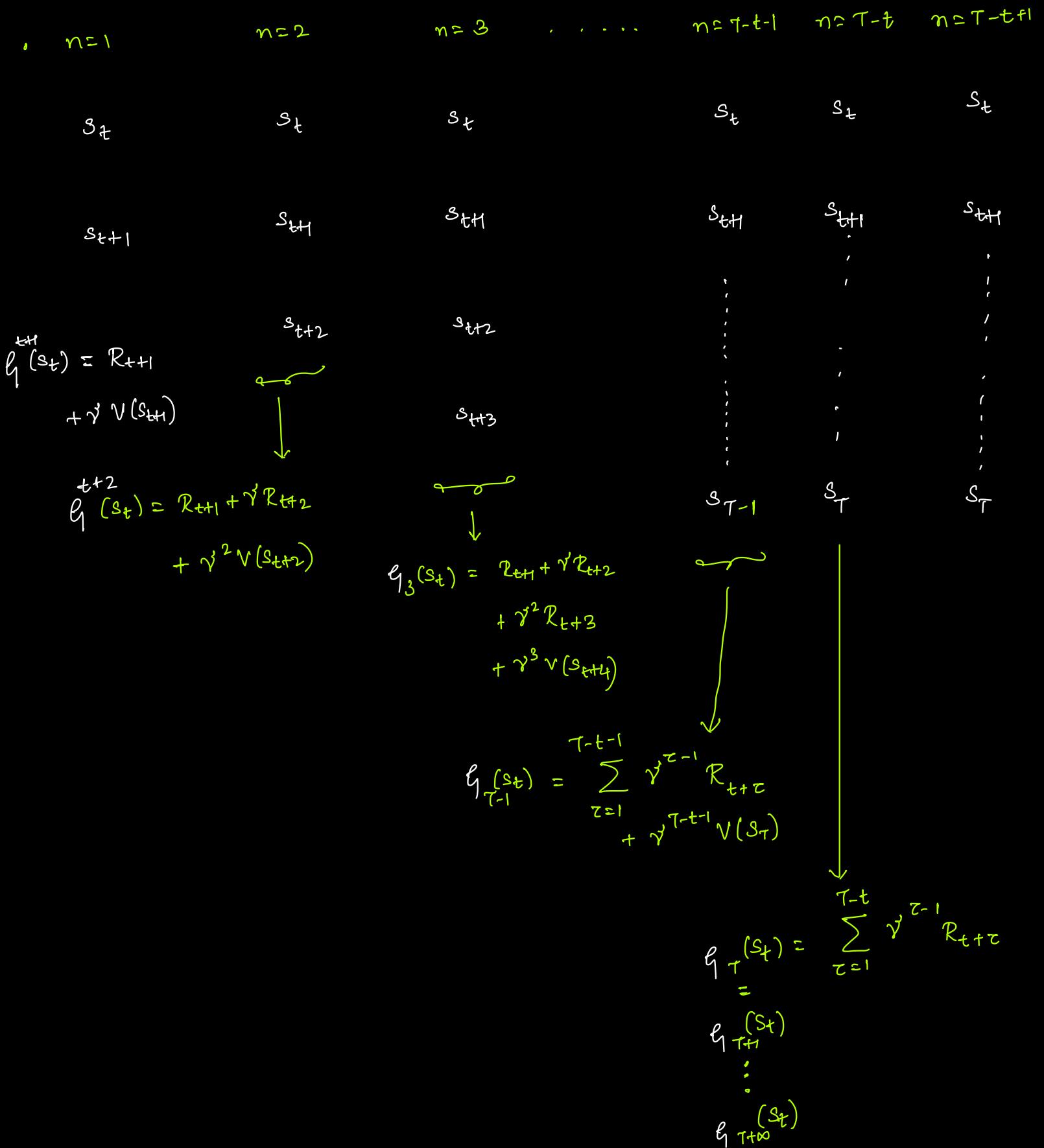
• n-step return:

$$V(s_t) = V(s_t) + \alpha \left( \sum_{\tau=t}^{t-n+1} \gamma^{t-\tau} \cdot R_{\tau+1} + \gamma^n V(s_{t+n}) - V(s_{t+n}) \right)$$

$\underbrace{\delta_n^{(s_t)} \text{ different } n\text{-step}}$

• Now consider a way to combine different n-step returns.

$$\Rightarrow V(s_t) = V(s_t) + \alpha \left[ \sum_{n=1}^{\infty} w_n \delta_n^{(s_t)} \right] \quad \textcircled{n} \rightarrow$$



$$\text{• } w \in \mathbb{R}^T, \quad \sum_{n=1}^{\infty} w_n = 1$$

$$\text{, let } w_n = \lambda^n \Rightarrow \sum_{n=1}^{\infty} w_n = \frac{1}{1-\lambda}$$

$$\Rightarrow (1-\lambda) \sum_{n=1}^{\infty} \lambda^n = 1$$

$$\therefore w_n = (1-\lambda) \lambda^n$$

$$\Rightarrow \underbrace{L_t}_{\text{Weighted return}} = (1-\lambda) \sum_{t=1}^{\infty} \lambda^{t-1} e_t^{t+n}$$

Weighted return

$$\Rightarrow L_t = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} e_t^{t+n} + (1-\lambda) \sum_{n=T-t}^{\infty} \lambda^{n-1} e_t^{t+n}$$

} for  $\lambda < 1$

$$= (1-\lambda) \sum_{t=1}^{T-t-1} \lambda^{n-1} e_t^{t+n} + \lambda^{T-t-1} e_t^T$$

• For  $\lambda = 0$  :

$$L_T = TD(0)$$

• For  $\lambda \rightarrow 1$  :

$$L_T \rightarrow MC = e_t^T \quad \left. \begin{array}{l} \text{Every visit } MC \\ \because v(s_t) = v(s_t) + \alpha e_t^T \end{array} \right\}$$

$\Rightarrow$  if  $s_t = s$  and  $s_{t+k} = s$

$\Rightarrow v(s)$  gets updated twice

NOTE :

• Offline Updating :

- Updates the value function estimates at the end of the trajectory. (store the correction / TD errors and update at the end)

• Online Updating :

- At every time step we update the value function estimates.

To update the value estimate of a state we need to wait till the end of the trajectory. i.e. we are looking at the forward states to update the value function of a state. Hence we call this as forward view or  $\lambda$ -return algorithm.

$$\text{i.e } V_{t+1}(s_t) = V_t(s_t) + \alpha [L_t - V_t(s_t)]$$

$$\text{where, } L_t = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} e_t^{t+n}$$

practically tedious to implement!

# BACKWARD VIEW (OR) TD( $\lambda$ ) ALGORITHM :

- The update rule is,

$$V_{t+1}(s_t) = V_t(s_t) + \alpha_t \delta_t e_t(s_t)$$

where,  $e_t(s) = \begin{cases} \gamma^\lambda e_{t-1}(s) + 1 & ; s = s_t \\ \gamma^\lambda e_{t-1}(s) & ; s \neq s_t \end{cases}$

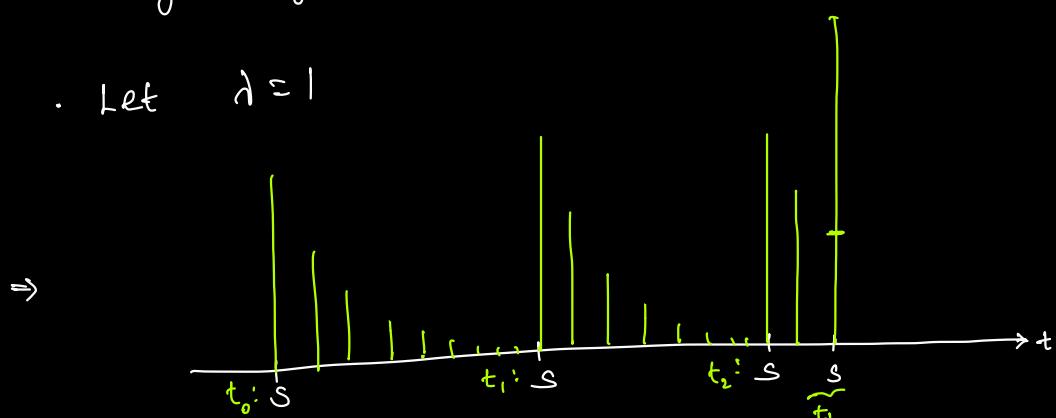
$$\delta_t = R_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)$$

- For offline TD( $\lambda$ ),

- Let  $\Delta(s) = \sum_{t=0}^T \delta_t e_t(s)$  ;  $\forall s$

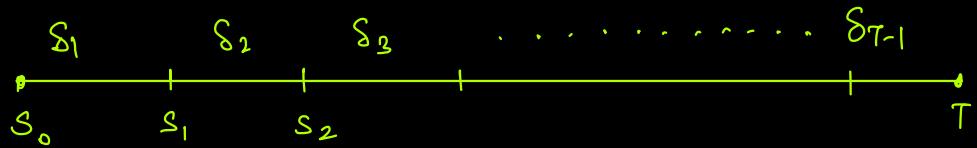
- How eligibility trace evolves?

- Let  $\lambda = 1$



Visit of state 's' at different time steps

The offline update for  $\lambda=1$  : (TD(1) algo<sup>m</sup>)



$$\Delta(s_0) = \delta_1 e_1(s_0) + \delta_2 e_2(s_0) + \delta_3 e_3(s_0) + \dots$$

$\downarrow$                                              $\downarrow$   
 $R_1 + \gamma V(s_1) - V(s_0)$                $R_2 + \gamma V(s_2) - V(s_1)$

$$\Rightarrow \Delta(s_0) = \left( R_1 + \gamma V(s_1) - V(s_0) \right) \underbrace{e_1(s_0)}_{\text{Assumption}} + \left( R_2 + \gamma V(s_2) - V(s_1) \right) \gamma + \dots$$

$$= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

$$+ \cancel{\gamma V(s_1)} - \cancel{\gamma V(s_1)} + \cancel{\gamma^2 V(s_2)} - \cancel{\gamma^2 V(s_2)} + \dots$$

$$- V(s_0) + \gamma^{T-1} V(s_T)$$

$$\boxed{\Delta(s_0) = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots + \gamma^{T-1} V(s_T) - V(s_0)}$$

16/3:

• Recap :

i) Forward view : ( $\lambda$ -return algorithm)

$s_t$

:

:

$s_{t+1}$

$s_t$

$s_{t+1}$

$s_{t+2}$

## EQUIVALENCE OF FORWARD VIEW AND BACKWARD VIEW :

The TD error from eligibility trace update rule is,

$$\Delta V_t^{TD}(s) = \alpha S_t e_t(s)$$

$$\text{where, } \delta_t = \gamma_{t+1} + \gamma^t V(s_{t+1}) - V(s_t)$$

$$e_t(s) = \gamma \lambda e_{t-1}(s) + \underbrace{\mathcal{I}_{ss_t}}_{\begin{cases} 1 & \text{if } s=s_t \\ 0 & \text{o.w} \end{cases}}$$

$$\begin{aligned} \Rightarrow e_t(s) &= \mathcal{I}_{ss_t} + \gamma \lambda (\gamma \lambda e_{t-2}(s) + \mathcal{I}_{ss_{t-1}}) \\ &= \mathcal{I}_{ss_t} + \gamma \lambda \mathcal{I}_{ss_{t-1}} + \gamma^2 \lambda^2 (e_{t-3}(s) + \mathcal{I}_{ss_{t-2}}) \\ &= \mathcal{I}_{ss_t} + \gamma \lambda \mathcal{I}_{ss_{t-1}} + \gamma^2 \lambda^2 \mathcal{I}_{ss_{t-2}} + \gamma^3 \lambda^3 \mathcal{I}_{ss_{t-3}} + \dots \end{aligned}$$

$$e_t(s) = \sum_{k=0}^t (\gamma \lambda)^{t-k} \mathcal{I}_{ss_k}$$

$\Rightarrow$  The sum of TD errors over the trajectory in backward view for a state 's' is,

$$\sum_{t=0}^{T-1} \Delta V_t^{TD}(s)$$

$$\Rightarrow \sum_{t=0}^{T-1} \Delta V_t^{\tau_D}(s) = \sum_{t=0}^{T-1} \alpha \delta_t e_t(s)$$

$$= \sum_{t=0}^{T-1} \alpha \delta_t \left( \sum_{k=0}^t (\gamma \lambda)^{t-k} I_{ss_k} \right)$$

- Changing 't' with 'k' and 'k' with 't'

$$= \sum_{k=0}^{T-1} \alpha \delta_k \left( \sum_{t=0}^k (\gamma \lambda)^{k-t} I_{ss_t} \right)$$

Expanding terms,

$$\begin{aligned}
&= \alpha \left[ \delta_0 \left( (\gamma \lambda)^0 I_{ss_0} \right) \right. \\
&\quad + \delta_1 \left( (\gamma \lambda) I_{ss_0} + (\gamma \lambda)^1 I_{ss_1} \right) \\
&\quad + \delta_2 \left( \gamma^2 \lambda^2 I_{ss_0} + \gamma \lambda I_{ss_1} + (\gamma \lambda)^2 I_{ss_2} \right) \\
&\quad + \delta_3 \left( \gamma^3 \lambda^3 I_{ss_0} + \gamma^2 \lambda^2 I_{ss_1} + \gamma \lambda I_{ss_2} + (\gamma \lambda)^3 I_{ss_3} \right) \\
&\quad \left. + \dots \right]
\end{aligned}$$

$$= \alpha \left( I_{SS_0} \sum_{k=0}^{T-1} (\gamma \lambda)^k \delta_k + I_{SS_1} \sum_{k=1}^{T-1} (\gamma \lambda)^{k-1} \delta_k + \dots \right)$$

$$= \alpha \sum_{t=0}^{T-1} I_{SS_t} \sum_{k=t}^{T-1} (\gamma \lambda)^{k-t} \delta_k$$

$$\Rightarrow \sum_{t=0}^{T-1} V_t^{TD}(s) = \sum_{t=0}^{T-1} \alpha I_{SS_t} \sum_{k=t}^{T-1} (\gamma \lambda)^{k-t} \delta_k$$

- ①

From the forward view the update rule is,

$$\Delta \hat{V}_t^\lambda(s_t) = \alpha \left[ \underbrace{L_t^\lambda - V_t(s_t)}_{\text{weighted sum of n-steps returns}} \right]$$

$$\Rightarrow \frac{1}{\alpha} \Delta \hat{V}_t^\lambda(s_t) = L_t^\lambda - V_t(s_t)$$

$$= -V_t(s_t)$$

$$+ (1-\lambda) \left[ R_{t+1} + \gamma V_t(s_{t+1}) \right]$$

$$+ (1-\lambda)\lambda \left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 V_t(s_{t+1}) \right]$$

⋮

$$= -V_t(s_t) (1-\lambda) \sum_{k=t}^{\infty} \lambda^{k-t} R_{t+1} + (1-\lambda)\lambda \gamma \sum_{k=t}^{\infty} \lambda^{k-t} \cdot R_{t+2} + \dots$$

$$\underbrace{\frac{1}{(1-\lambda)}}$$

$$+ \gamma V_t(s_{t+1}) + \gamma^2 \lambda V_t(s_{t+2}) + \dots$$

$$- \lambda \gamma V_t(s_{t+1}) - \gamma^2 \lambda^2 V_t(s_{t+2}) + \dots$$

$$= R_{t+1} + \gamma \lambda R_{t+2} + \dots$$

$$+ \gamma V_t(s_{t+1}) - V_t(s_t)$$

$$+ \gamma^\lambda \left[ \gamma^\lambda v_t(s_{t+2}) - v_t(s_{t+1}) \right]$$

$$+ \gamma^2 \lambda^2 \left[ \gamma^\lambda v_t(s_{t+3}) - v_t(s_{t+2}) \right]$$

+ ...

$\delta_t$

$$= (\gamma^\lambda)^\circ \left[ R_{t+1} + \gamma^\lambda v_t(s_{t+1}) - v_t(s_t) \right]$$

$$+ (\gamma^\lambda) \left[ R_{t+2} + \gamma^\lambda v_t(s_{t+2}) - v_t(s_{t+1}) \right]$$

$\delta_{t+1}$

+ ...

$$= \sum_{k=t}^{\infty} (\gamma^\lambda)^{k-t} \delta_k \quad \left. \begin{array}{l} \text{For offline update} \\ \text{equality holds} \end{array} \right\}$$

$$\Rightarrow \Delta v_t^\lambda(s_t) = \alpha \sum_{k=t}^{\infty} (\gamma^\lambda)^{k-t} \delta_k$$

For a state  $s$  the update over a trajectory under forward view is,

$$\sum_{t=0}^{T-1} \Delta v_t^\lambda(s) I_{ss_t} = \sum_{t=0}^{T-1} \alpha \sum_{k=t}^{\infty} (\gamma^\lambda)^{k-t} \delta_k \cdot I_{ss_t}$$

Since the rewards are zero  $\forall k > T$

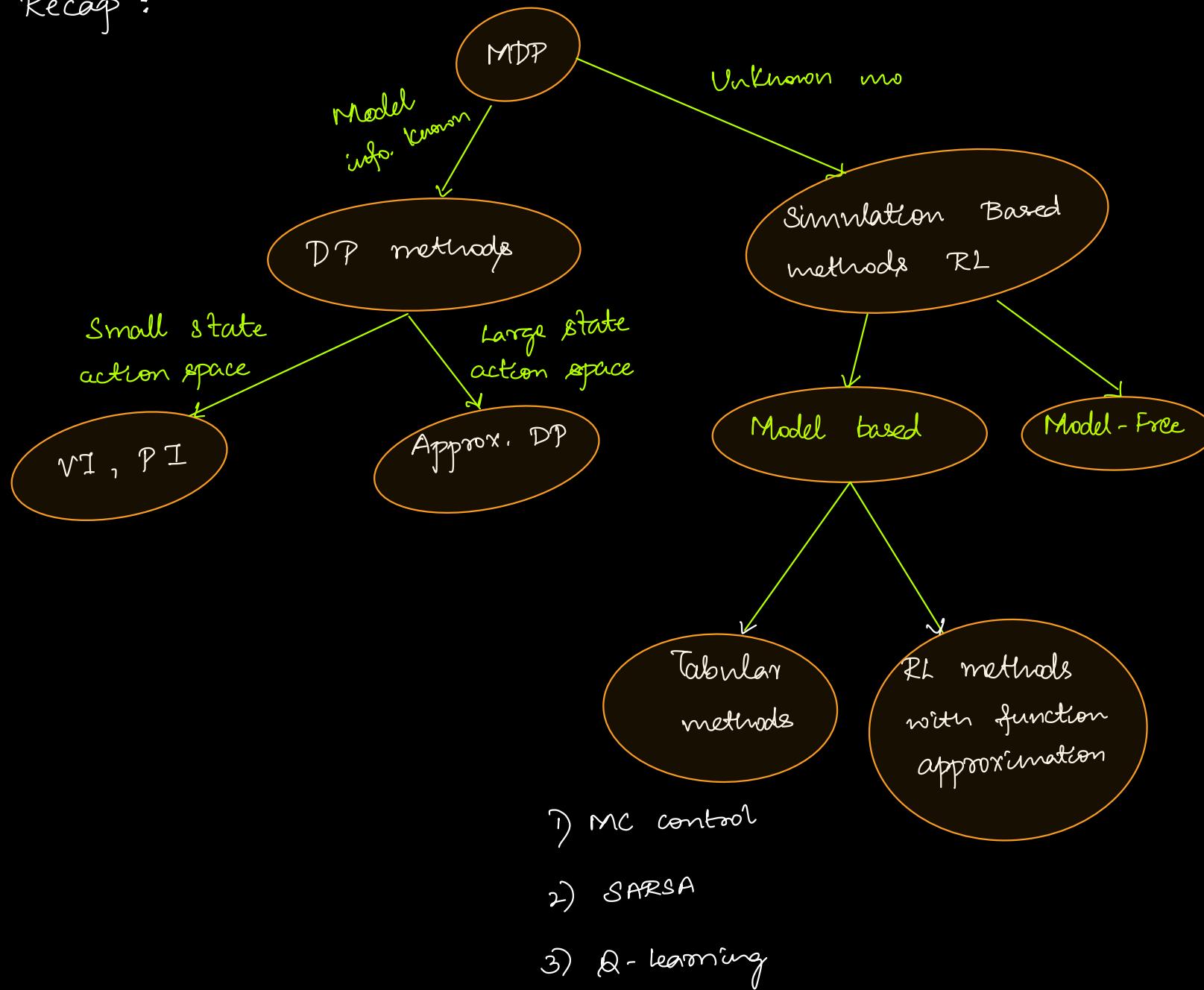
$$\Rightarrow \sum_{t=0}^{T-1} \Delta v_t^\lambda(s) I_{ss_t} = \sum_{t=0}^{T-1} \alpha \sum_{k=t}^{T-1} (\gamma^\lambda)^{k-t} \delta_k I_{ss_t} - \textcircled{2}$$

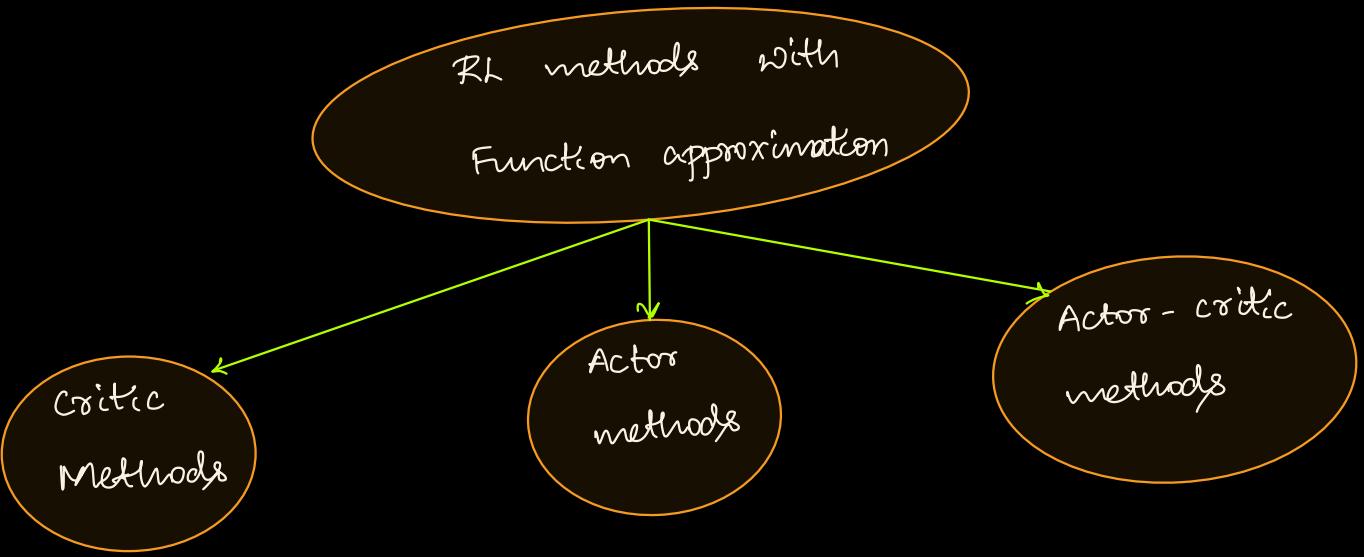
From ① and ②

$$\sum_{t=0}^{T-1} \Delta v_t^{\text{TD}}(s) = \sum_{t=0}^{T-1} \Delta v_t^\lambda(s) I_{ss_t}; \quad \forall s \in S$$

18/3 :

Recap :





• What functions we can approximate ?

- 1) Value functions
- 2) Policy

VALUE - BASED FUNCTION APPROXIMATION :

• Why approximation ?

• Due to large-state space we can't store  $v(s)$  &  $s$  and computing  $v(s)$  &  $s$  is computationally expensive.

• How are we storing for continuous state space .

Ex:

$$T \xrightarrow{-1} A \xrightarrow{0} B \xrightarrow{0} C \xrightarrow{0} D \xrightarrow{0} E \xrightarrow{1} T$$

- One trajectory is  $CDEDT$

$$v(E) = \gamma_2 \cdot 1 + \frac{1}{2} v(D)$$

$$v(D) = \gamma_2 v(E) + \gamma_2 v(C)$$

$$v(C) = \gamma_2 v(B) + \gamma_2 v(D)$$

$$v(B) = \gamma_2 v(C) + \gamma_2 v(A)$$

$$v(A) = \gamma_2 \cdot 0 + \gamma_2 v(B)$$

$$\Rightarrow V_B = 2V_A$$

- A linear approximation for a fixed policy corresponds to features we considered

$$\hat{v}_\mu(s, w) = \sum_{i=1}^k w_i \phi_i(s, w)$$

Basis function  
(can be polynomial | RBF)

- Our goal is  $\min \sum_{s \in S} (V_M(s) - \hat{v}_\mu(s))^2$   
gives equal importance  
 $\forall s \in S$

We give more weightage to the states that are visited often.

$$\Rightarrow \text{goal : } \min \sum_{s \in S} d_{\mu}(s) \left( V_{\mu}(s) - \hat{V}_{\mu}(s) \right)^2$$

↓  
Stationary dist. of state  $s$

For linear approximator,

$$J(\omega) = \sum_{s \in S} d_{\mu}(s) \left( V_{\mu}(s) - \phi(s)^T \omega \right)^2$$

where,  $\phi(s) = (\phi_1(s), \phi_2(s), \dots, \phi_k(s))$

$$\text{goal : } \min_{\omega} J(\omega)$$

Using gradient descent :

$$- 2 \sum_{s \in S} d_{\mu}(s) (V_{\mu}(s) - \phi(s)^T \omega) \phi(s)$$

$$\Rightarrow \omega = \omega + \alpha \cdot 2 \cdot \sum_{s \in S} d_{\mu}(s) \underbrace{\left( V_{\mu}(s) - \phi(s)^T \omega \right)}_{\hat{V}_{\mu}(s)} \phi(s)$$

- The stochastic update rule is,

$$w_{n+1} = w_n + \alpha \left[ \underbrace{\left( V_\mu(s) - \hat{V}(s, w) \right)}_{\text{Unknown}} \phi(s) \right]$$

- MC function Approximation :

$$w_{n+1} = w_n + \alpha \left[ e_t(s) - \hat{V}(s, w) \right] \phi(s)$$

- TD(0) Function Approximation :

$$w_{n+1} = w_n + \alpha \left[ R_{t+1} + \gamma \hat{V}(s', w) - \hat{V}(s, w) \right] \phi(s)$$

- TD( $\lambda$ ) Function Approximation :

$$w_{n+1} = w_n + \alpha \left[ R_{t+1} + \gamma \hat{V}(s', w) - \hat{V}(s, w) \right] e_t$$

where,  $e_t = \gamma^\lambda e_t + \phi(s_t)$