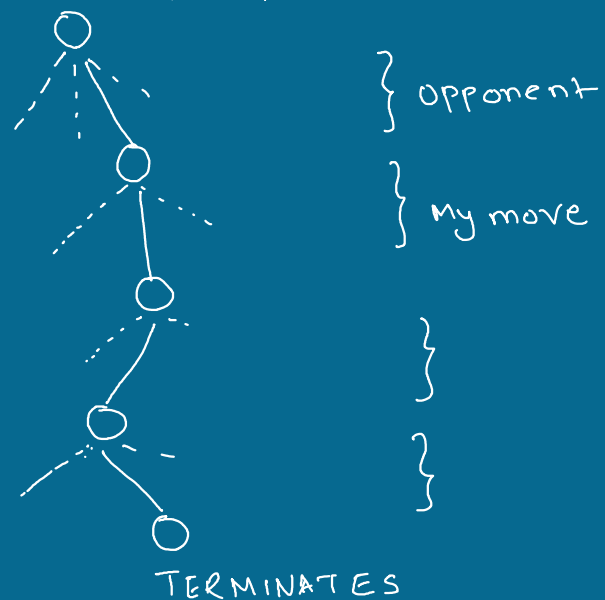


Value Backup



- Fully informed
- Partially informed

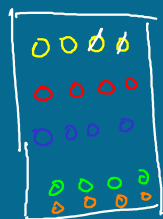
History

MENACE
(1961-62)

trial-and-error learning

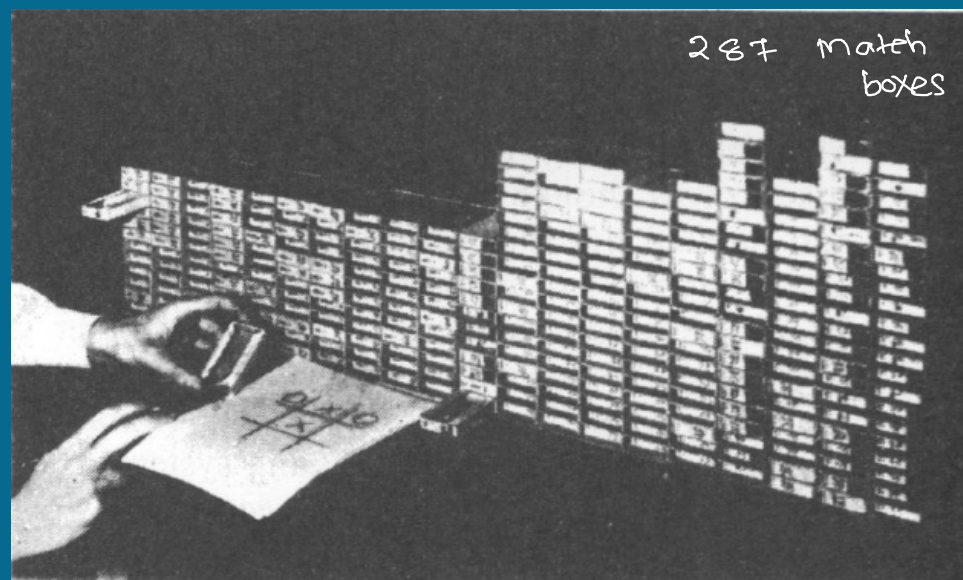
(Donald Michie)

x	②	③
0	x	④
①	0	⑤



Menace — 'O'

I am — 'x'



Matchbox Educable Noughts & Crosses Engine.

N-Arm bandit :

- Stationary distribution
- "Explore-Exploit" time period

Choices (Actions)

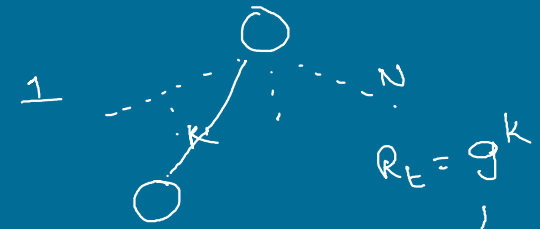
$R_t \leftarrow$ Reward from Environment

t 1 2 3 ... N

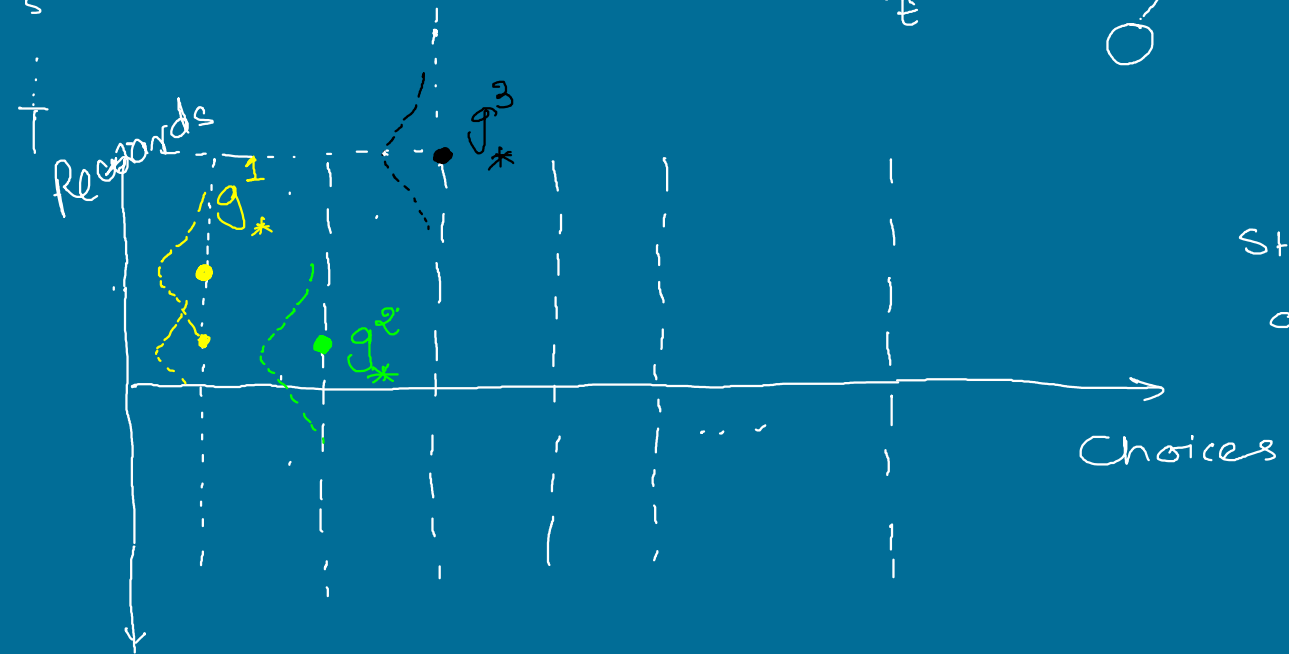
1
2
3 \equiv Expected total reward
4 $E[R_1 + R_2 + \dots + R_T]$
5

$t=1$

$a_{t=1} = k$



Stationary distribution



$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

A simple Bandit Algorithm :

- stationary
- Non-stationary



- $Q_{n+1} = Q_n + \alpha [R_n - Q_n]$

↑
fixed

(weighted average)

$$Q_{n+1} = \underbrace{(1-\alpha)^n}_{(1-\alpha)^n + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} = 1} Q_1 + \sum_{i=1}^n \underbrace{\alpha(1-\alpha)^{n-i}} R_i$$

- Initialization

$$Q(a) \leftarrow 0 \quad \forall a$$

$$N(a) \leftarrow 0 \quad \forall a$$

Loop forever

$$A \leftarrow \begin{cases} \arg\max_a Q(a) & \sim 1-\epsilon \\ \text{'a' random} & \sim \epsilon \end{cases}$$

$$R \leftarrow \text{Bandit}(A)$$

$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{n} [R - Q(A)]$$

$$\sum_{n=1}^{\infty} \alpha_n(a) = \infty$$

and

$$\sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$$

Reading Assignment:

2.6 Opt Init Values

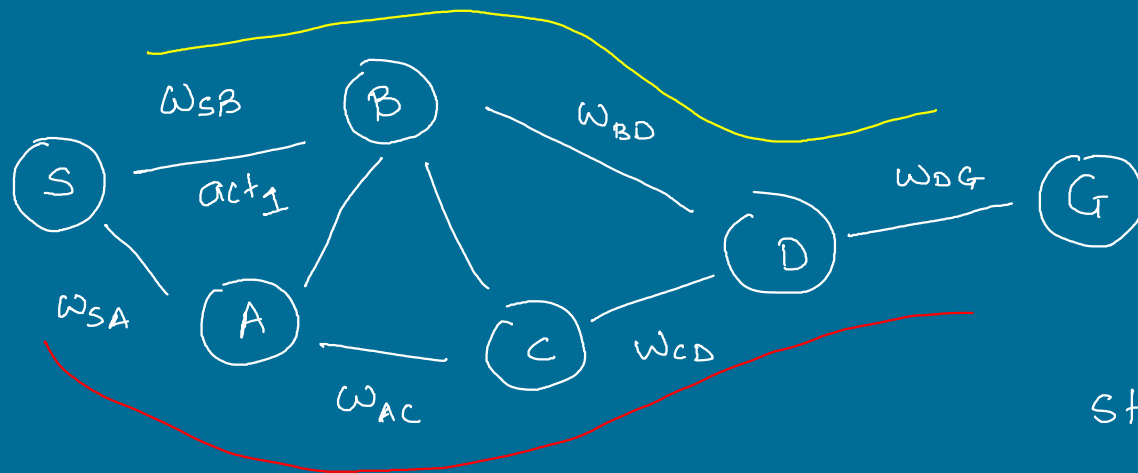
2.7 UCB Action Selection

2.8 Gradient Bandit Algorithm

Markov Decision Process:

(sequential)

Typical Search

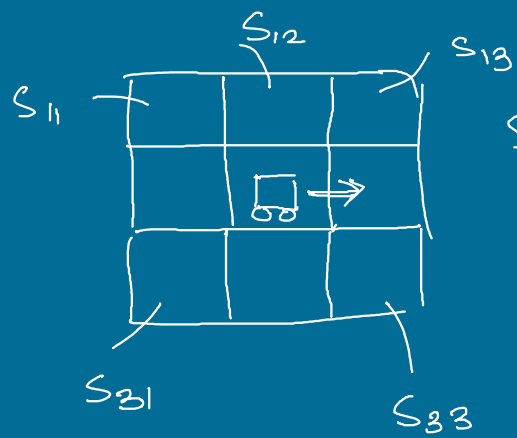


action (state)

succ (state, action)

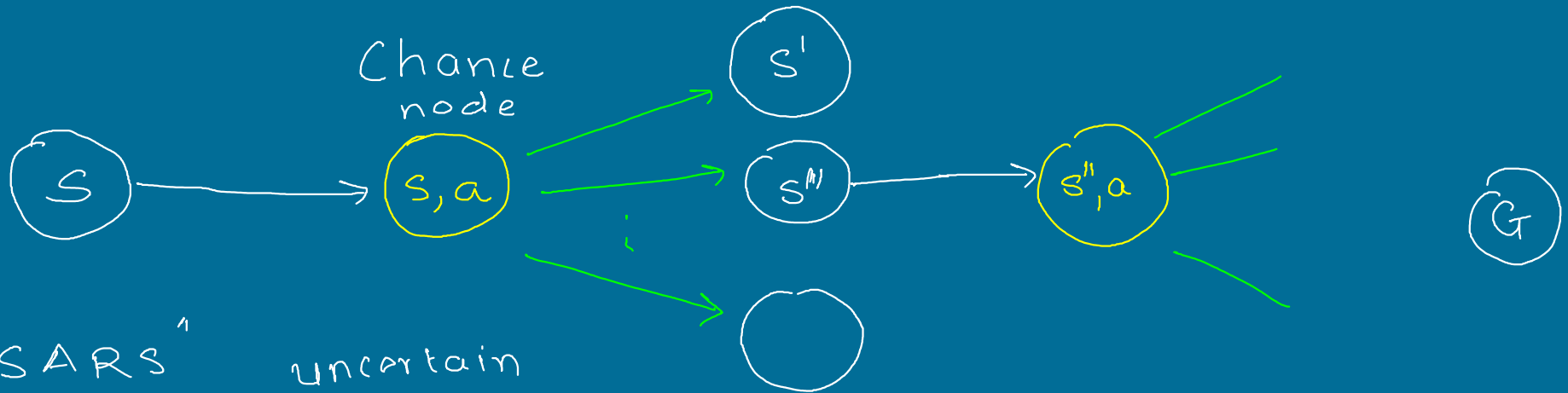
1 certain 2
 state, act, weight, state, act, w

□ ... , state = G



$\text{succ}(s_{22}, \rightarrow) \rightsquigarrow$

	Probability
s_{23}	0.5
s_{13}	0.05
s_{33}	0.05
s_{12}	0.2
s_{32}	0.2



"SARS"

uncertain

$s_{t=0}, A_{t=0}, R_{t=1}, s_{t=1}, A_{t=1}, R_{t=2}, \dots$

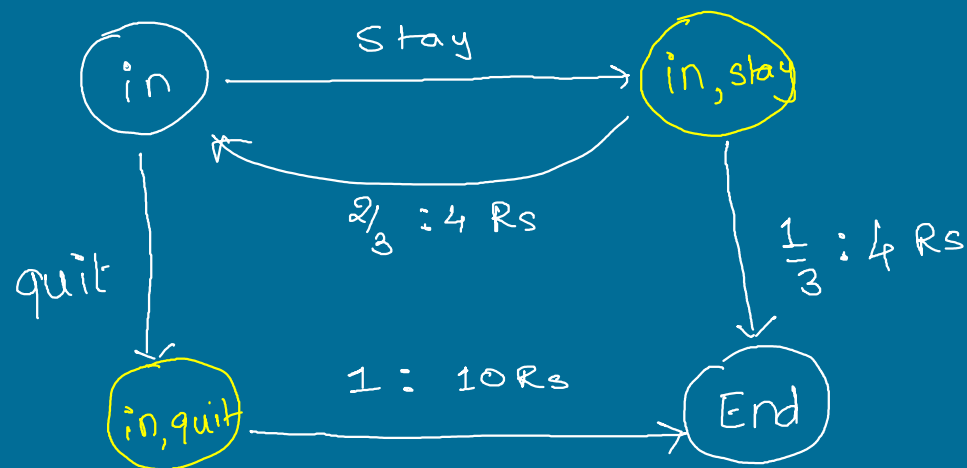
non-deterministic

uncertain

Example: Dice Game

- You choose to stay / quit
- If you quit, you get 10 Rs and game ends
- if you stay, you get 4 Rs and then environment rolls a dice

Expected Total Reward + If the dice results in 2/3, game ends
+ Otherwise, continue



Policy

$\pi(\text{state}) = \text{action}$

$\pi(\text{in}) = \text{stay}$

$\pi(\text{end}) \neq ?$

$\pi(\text{in}) = \text{quit}$

10 Rs.

Run 1: $S_0 = \text{in}$, $A_0 = \text{stay}$, $R_1 = 4R_s$, $S_1 = \text{in}$, $A_1 = \text{stay}$, $R_2 = 4R_s$,
 $S_2 = \text{END}$.

\tilde{G}_0^1 8

Run 2: $S_0 = \text{in}$, $[A_0 = \text{stay}, R_1 = 4R_s, S_1 = \text{in}]$, $[A_1 = \text{stay}, R_2 = 4R_s,$
 $S_2 = \text{in}]$, $[A_2 = \text{stay}, R_3 = 4R_s, S_3 = \text{in}]$, $[A_3 = \text{stay}, R_4 = 4R_s,$
 $S_4 = \text{END}]$.

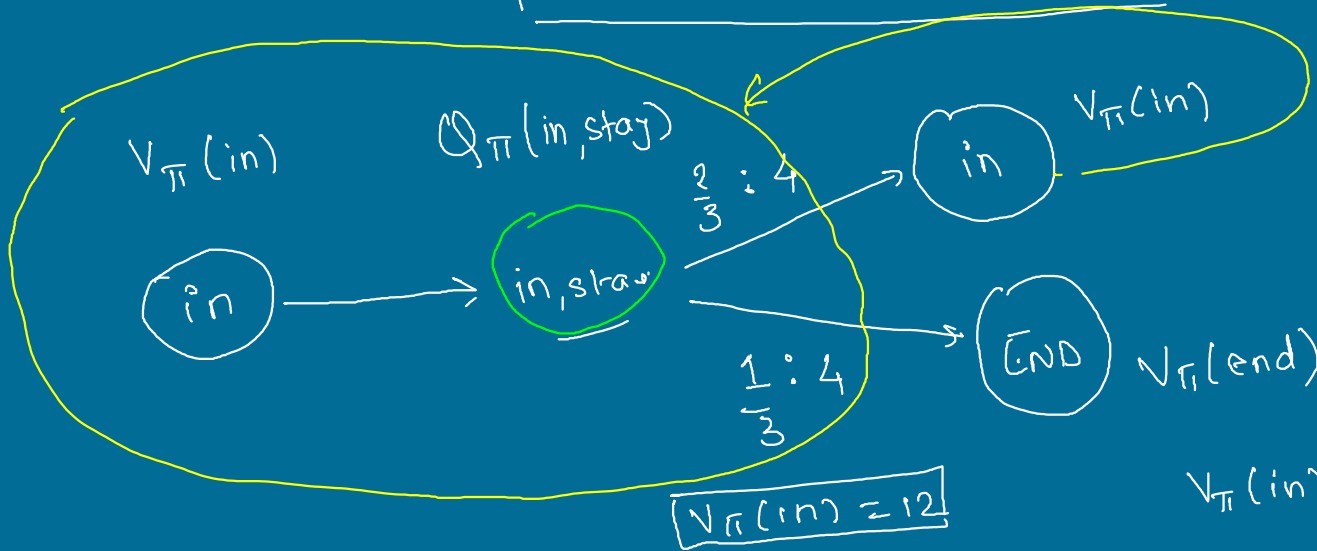
\tilde{G}_0^2 16

Run N

12

$$\frac{1000 \times 10}{1000 \times 12}$$

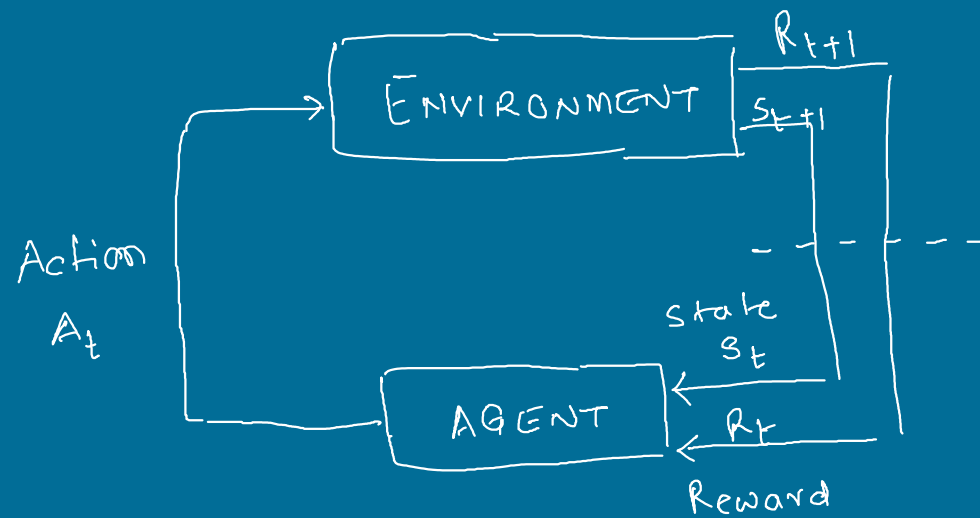
Avg Total Reward \sim Expected Total Reward \Rightarrow Expected Return $S = \text{in}$
 $\pi(\text{in}) = \text{stay}$



$$V_\pi(\text{in}) = \frac{2}{3} (4 + V_\pi(\text{in})) + \frac{1}{3} (4 + V_\pi(\text{end}))$$

$$V_\pi(\text{in}) = 4 + \frac{2}{3} V_\pi(\text{in})$$

Agent Environment Interface



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

$$s_0 \quad A_0 \quad \underbrace{R_1 \quad s_1 \quad A_1}_{t \quad t+1} \quad \underbrace{R_2 \quad s_2 \quad A_2 \quad \dots}_{\text{Trajectory}}$$

$$s_t \in S, \quad A_t \in A(s_t), \quad R_t \in \mathbb{R} \subseteq \mathbb{R}$$

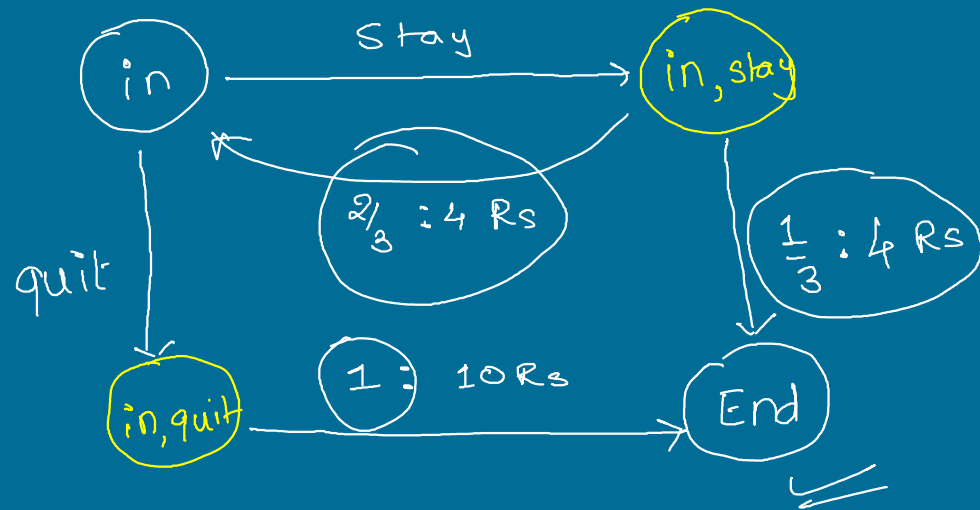
Finite MDP : $|S| < \infty, |A(s)| < \infty, |R| < \infty$

$$p(s_{t+1} = s', R_{t+1} = r \mid s_t = s, A_t = a)$$

$$p : S \times \mathbb{R} \times S \times A \rightarrow [0, 1]; \quad \sum_{s' \in S} \sum_{r \in \mathbb{R}} p(s', r \mid s, a) = 1$$

$$s \in S, a \in A(s).$$

$$V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a)$$



$\Gamma = 1$
deterministic
 $\pi(\text{in}) = \text{stay}$
 $\pi(\text{end}) = \{ \}$

$$V_{\pi}(\text{in}) = 12$$

$$V_{\pi}(\text{End}) = 0$$

$\Gamma = 1$
nondeterministic
 $\pi(\text{stay}|\text{in}) = 0.5$
 $\pi(\text{quit}|\text{in}) = 0.5$
 $V_{\pi}(\text{in}) = ?$
 $V_{\pi}(\text{end}) = 0$

$$V_{\pi}(\text{in}) = \left\{ \overset{1/2}{\pi(\text{stay}|\text{in})} \left(\overset{2/3}{p(\text{in}, 4|\text{in}, \text{stay})} \cdot [4 + V_{\pi}(\text{in})] + \overset{1/3}{p(\text{end}, 4|\text{in}, \text{stay})} \cdot [4 + V_{\pi}(\text{end})] \right) + \right.$$

$$\left. \overset{1/2}{\pi(\text{quit}|\text{in})} \left(\overset{1}{p(\text{end}, 10|\text{in}, \text{quit})} \cdot [10 + V_{\pi}(\text{end})] \right) \right\}$$

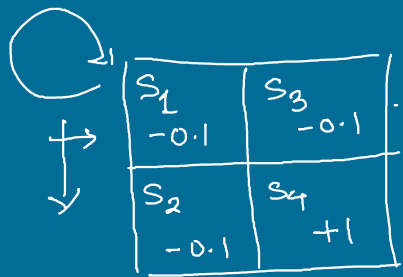
$$V_{\pi}(\text{in}) = \textcircled{\frac{1}{2}} \left[\textcircled{\frac{2}{3}} \cdot [4 + V_{\pi}(\text{in})] + \frac{1}{3} \cdot 4 \right] + \textcircled{\frac{1}{2}} \left[\textcircled{1} \cdot 10 \right] \neq$$

(A) (B)

$$2V_{\pi}(in) = \frac{2}{3} V_{\pi}(in) + 14$$

$$V_{\pi}(in) = \frac{3}{4} \times 14 = 10.5$$

$$\begin{matrix} < 12 \\ > 12 \\ \approx 12 \end{matrix}$$

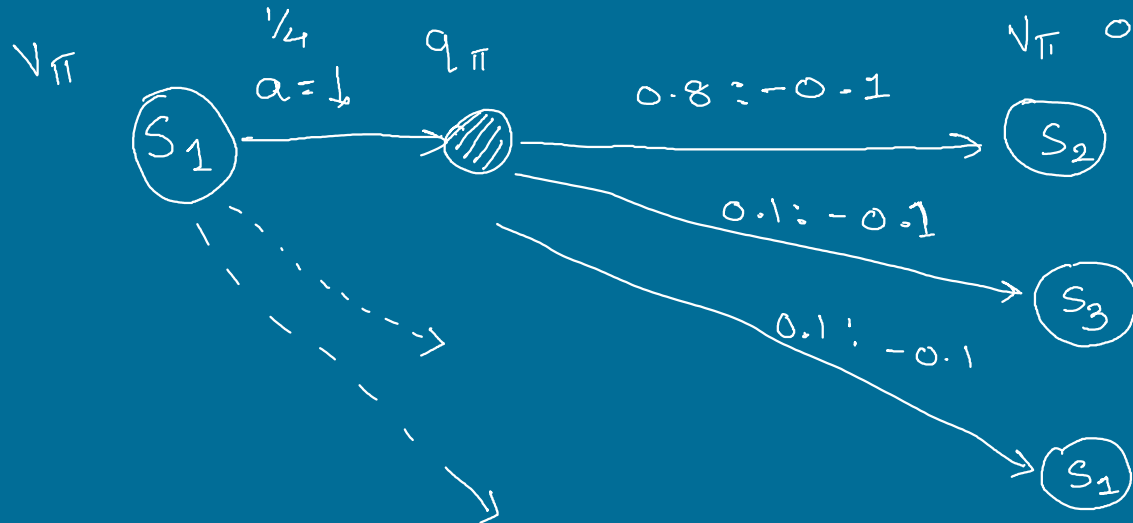
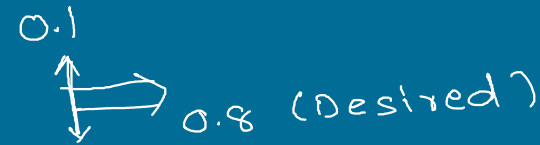


Grid World Example:

$$\pi(a|s) = 1/4$$

Policy $a \in \{\uparrow, \downarrow, \rightarrow, \leftarrow\}$

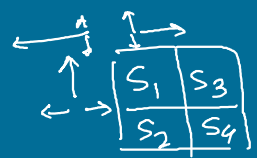
Environment is stochastic:



Value Function (state)

$$V_{\pi}(s) = \begin{cases} +1 & \text{if } s = s_4 \\ \sum_a \pi(a|s) q_{\pi}(s, a) & \end{cases}$$

(s is not a terminal state)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$


$$V_{\pi}(s_1) = \frac{\pi(\uparrow | s_1)}{1/4} \left[\overbrace{0.8 (-0.1 + \gamma V_{\pi}(s_1)) + 0.1 (-0.1 + \gamma V_{\pi}(s_1))} + 0.1 (-0.1 + \gamma V_{\pi}(s_3)) \right]$$

$$+ \frac{\pi(\downarrow | s_1)}{1/4} \left[0.8 (-0.1 + \gamma V_{\pi}(s_2)) + 0.1 (-0.1 + \gamma V_{\pi}(s_3)) + 0.1 (-0.1 + \gamma V_{\pi}(s_1)) \right]$$

$$+ \frac{\pi(\rightarrow | s_1)}{1/4} \left[0.8 (-0.1 + \gamma V_{\pi}(s_3)) + 0.1 (-0.1 + \gamma V_{\pi}(s_1)) + 0.1 (-0.1 + \gamma V_{\pi}(s_2)) \right]$$

$$+ \frac{\pi(\leftarrow | s_2)}{1/4} \left[0.9 (-0.1 + \gamma V_{\pi}(s_1)) + 0.1 (-0.1 + \gamma V_{\pi}(s_2)) \right]$$

$$V_{\pi}(s_1) = \frac{1}{4} V_{\pi}(s_1) + \frac{1}{4} V_{\pi}(s_2) + \frac{1}{4} V_{\pi}(s_3) + \frac{1}{4} V_{\pi}(s_4) + C_1 \quad (1)$$

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$(V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')])$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

(I) Two arrays for maintaining $v^{\text{iter}}, v^{\text{iter}+1}$

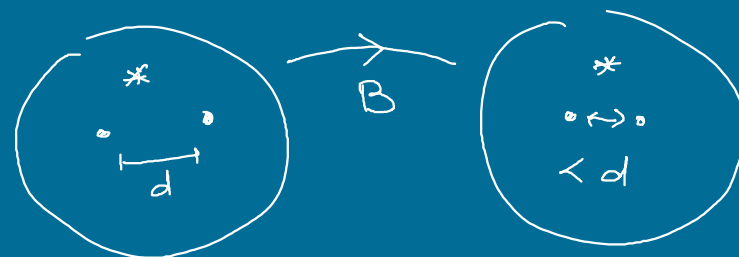
(II) In-place

Q. Convergence? "Contraction" and "Fixed point th^m"

$$V^0 = V^{\text{init}}$$

$$V^{i+1} = \underbrace{B}_{\substack{\uparrow \\ \text{contraction}}} V^i = (B^i) V^{\text{init}}$$

$$B : \mathbb{R}^4 \rightarrow \mathbb{R}^4$$



$$B \underline{*} = \underline{*}$$

08.04.2021

Summary

- Finite MDP
- Value functions $\leftarrow \begin{matrix} \text{state value} \\ \text{state action} \end{matrix}$

$$V_{\pi}(s)$$

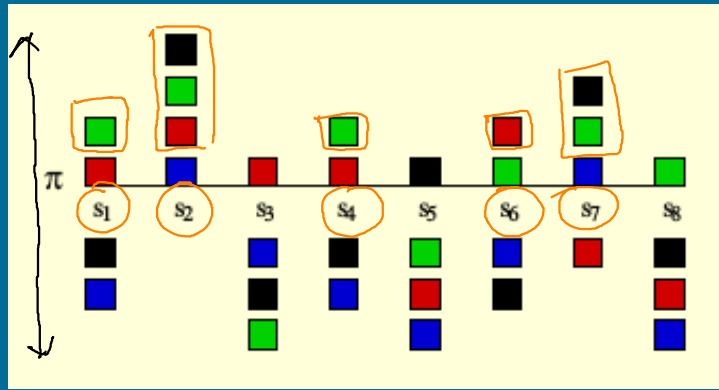
$$Q_{\pi}(s, a)$$

$$V_{\pi}(s)$$



[Bellman Eqn.]

Policy Improvement: Given a policy π

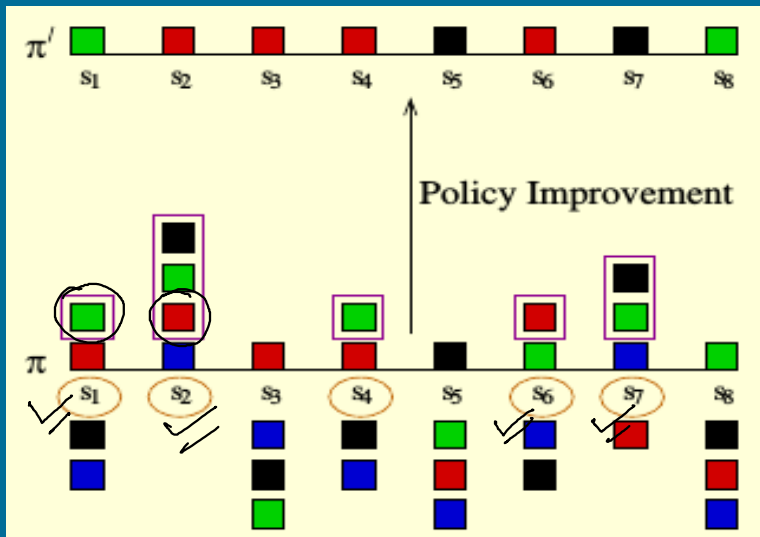


s_1 s_8

Improvable states s_1, s_2, s_4, s_6, s_7

Improvable actions $\{ [G], [B, G, R], [G], [R], [B, G] \}$ improving action

- Pick one or more improvable states and in them
- Switch to an arbitrary
- Let the resulting policy be π' .

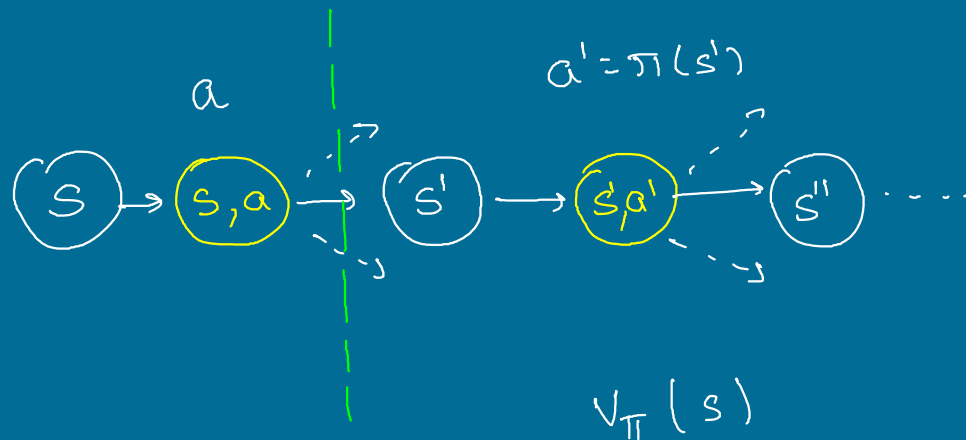


Policy Improvement Theorem:

(1) If π has no improvable states, then it is OPTIMAL else

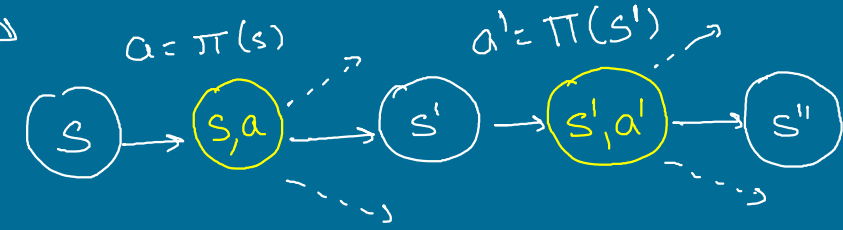
(2) If π' is obtainable as above then
 $\forall s \in S: V^{\pi'}(s) \geq V^{\pi}(s), \exists s \in S: V^{\pi'}(s) > V^{\pi}(s)$

$$Q_{\pi}(s, a) \quad \left\{ \begin{array}{l} Q_{\pi}(s, \pi(s)) = V_{\pi}(s) \end{array} \right.$$



$$Q_{\pi}(s, a) > Q_{\pi}(s, \pi(s))$$

π' is same as π everywhere except at s , $\pi'(s) = a$



$s \neq \text{terminal}$

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid s_t = s]$$

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

$$Q_{\pi'}(s, \pi'(s)) > Q_{\pi}(s, \pi(s))$$

$$a = \underset{\text{act}}{\operatorname{argmax}} Q_{\pi}(s, \text{act})$$

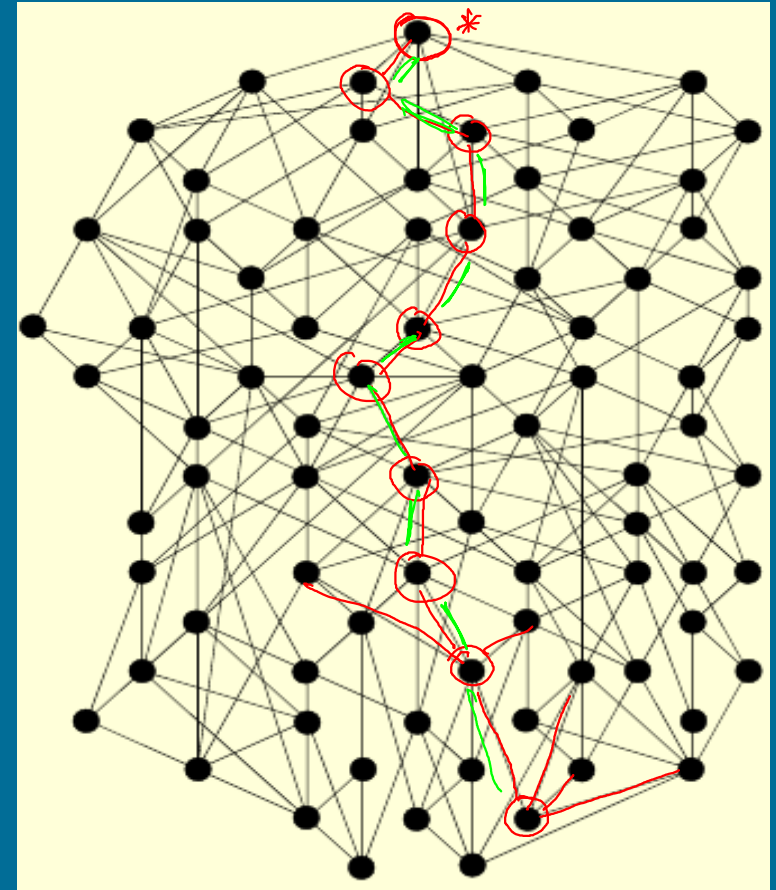
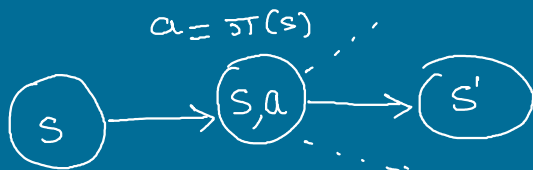
- $X: S \rightarrow \mathbb{R}$ and $Y: S \rightarrow \mathbb{R}$, $\boxed{X \geq Y}$ if $\forall s \in S; X(s) \geq Y(s)$.
- $X > Y$ if $X \geq Y$ and $\exists s \in S; X(s) > Y(s)$.

- For policies $\pi_1, \pi_2 \in \Pi$, $\pi_1 \geq \pi_2$ if $V^{\pi_1} \geq V^{\pi_2}$
and $\pi_1 > \pi_2$ if $V^{\pi_1} > V^{\pi_2}$

- Bellman Operator:

$$B^\pi: (S \rightarrow \mathbb{R}) \rightarrow (S \rightarrow \mathbb{R})$$

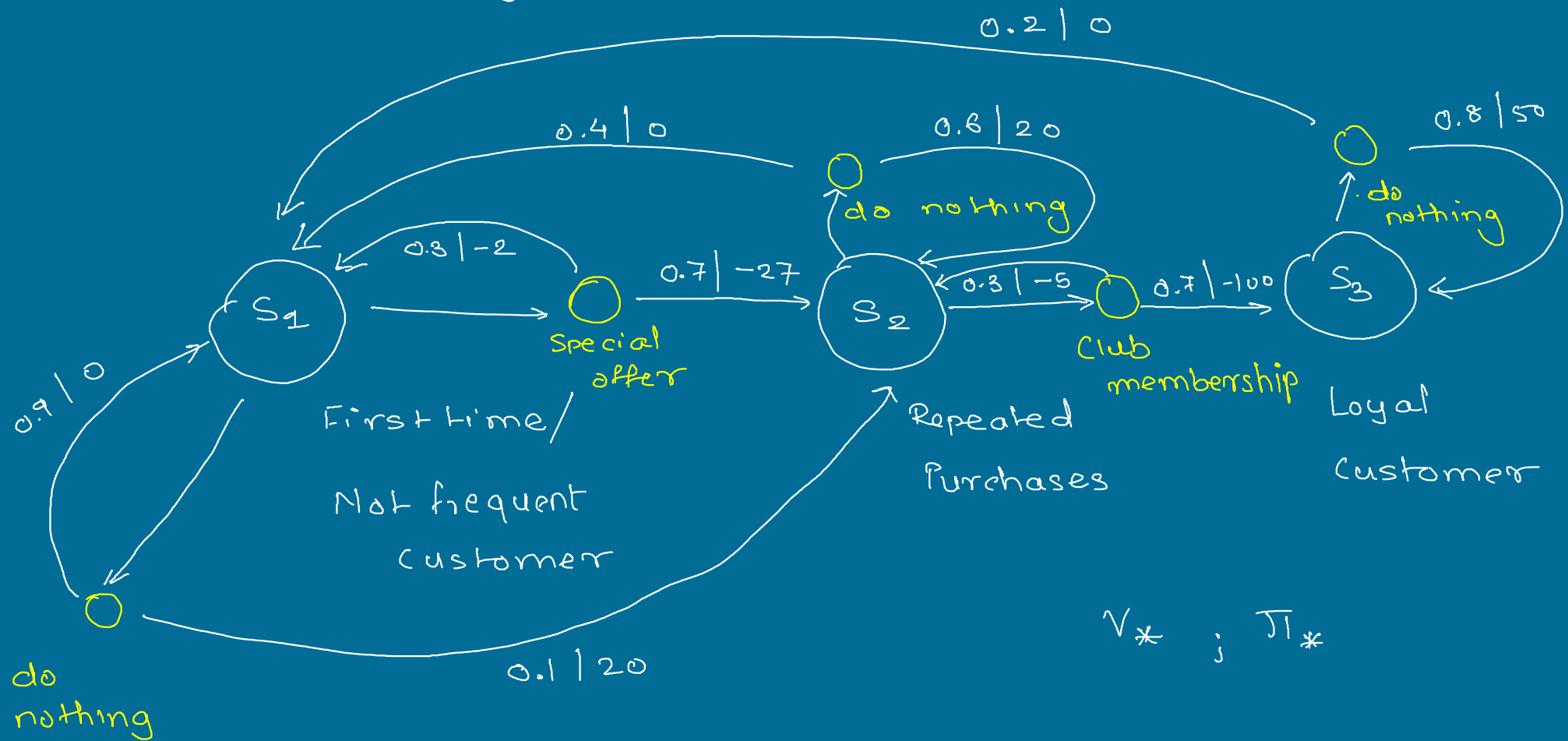
$$V_\pi(s) = \sum_{s', r} p(s', r | s, \pi(s)) [\gamma + \gamma V_\pi(s')]$$



"Policy PosET"

$$B^{(\pi)}(\underline{x(s)}) = \sum_{s', r} p(s', r | s, \pi(s)) [\gamma + \gamma x(s')]$$

Example: (Advertising Problem)



1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

$$\pi \rightarrow V^\pi$$

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

$$\pi \rightarrow \pi'$$

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Policy Iteration

$$\pi_1 \downarrow \text{Eval}$$

$$V^{\pi_1}$$

$$\downarrow \text{Imp}$$

$$\pi_2$$

$$\downarrow \text{Eval}$$

$$V^{\pi_2}$$

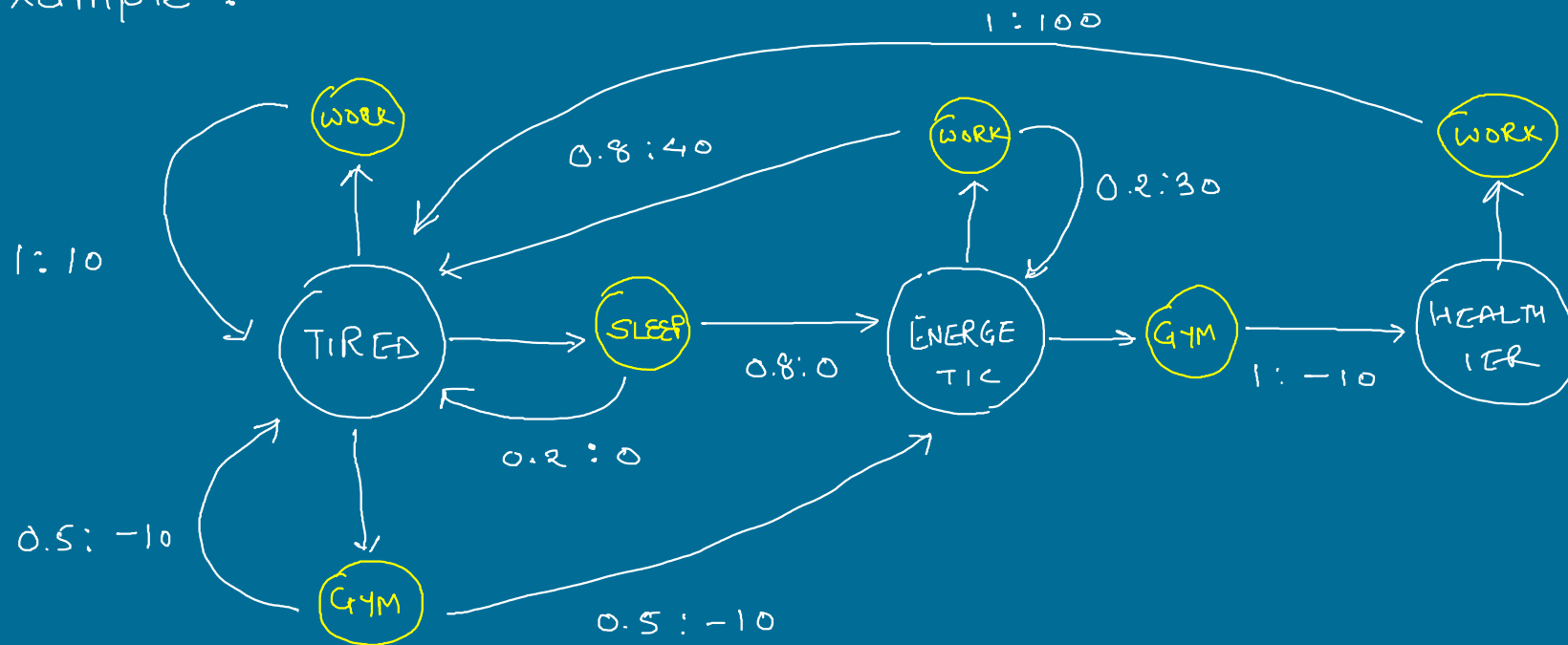
$$\downarrow \text{Imp}$$

$$\pi_3$$

$$\vdots$$

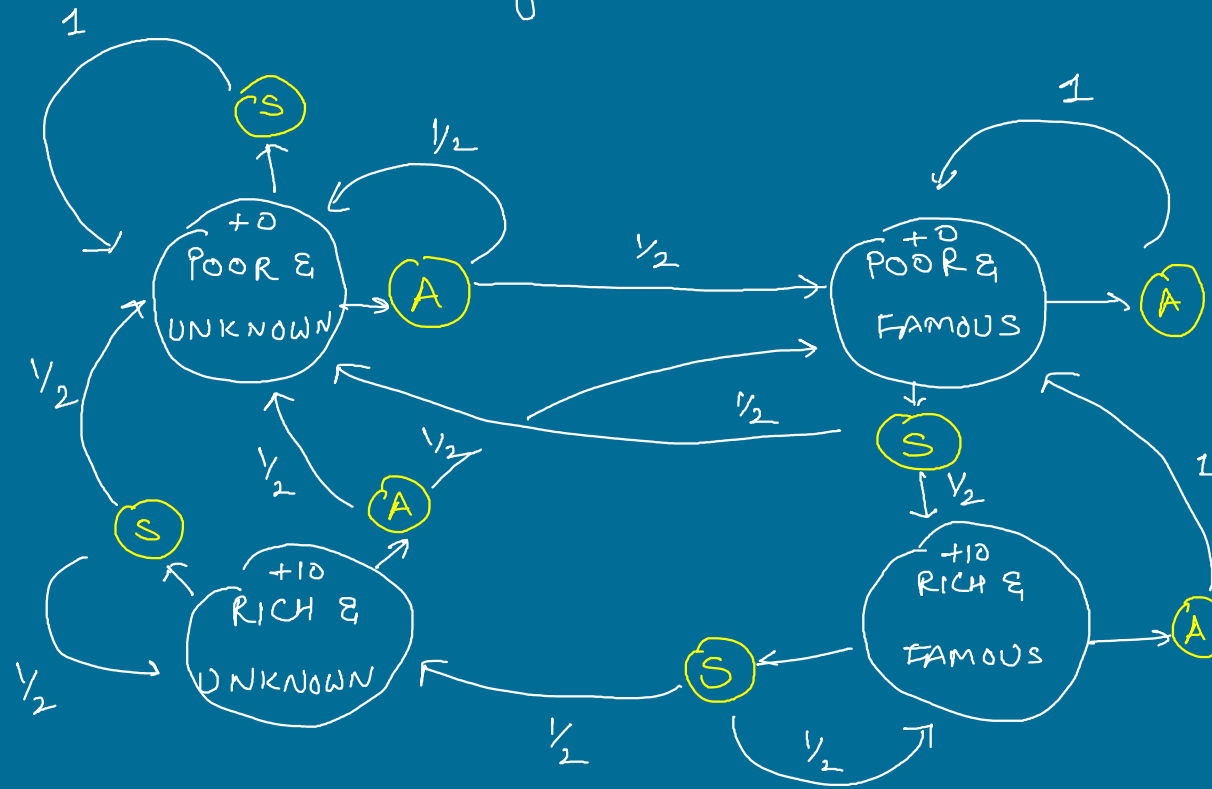
$$(V_*, \pi_*)$$

Example : (Actor's Life cycle)



Example:

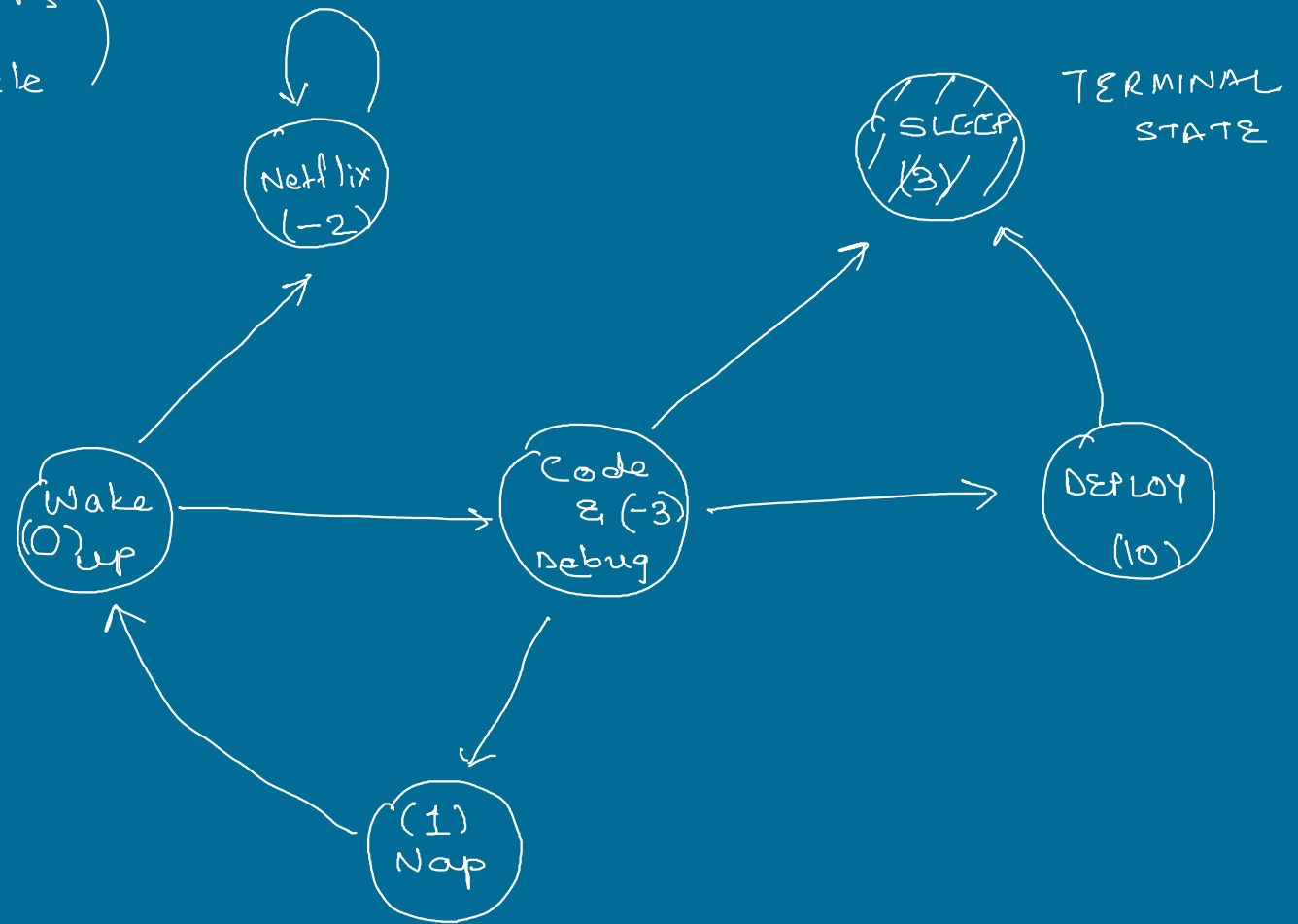
(Company Outreach)



S - Save

A - Advertise

Example : (Developer's Life Cycle)



Optimal Bellman Equations

$$V^*(s) = \max_a \left[\sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')] \right] \quad (\text{Value Iteration})$$

$$Q^*(s, a) =$$

[Exercise]

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

| $\Delta \leftarrow 0$

| Loop for each $s \in \mathcal{S}$:

| $v \leftarrow V(s)$

| $V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

| $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

Monte Carlo Methods:

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

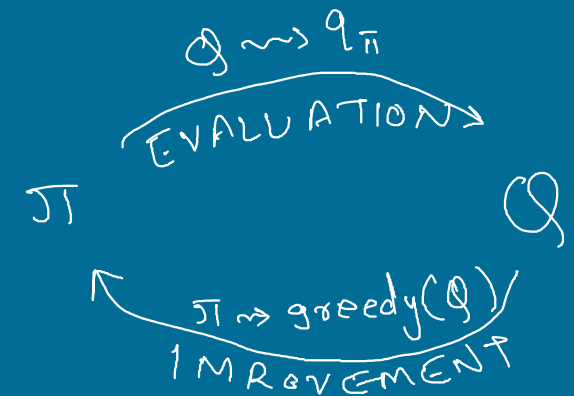
Append G to $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

- Model is not known

Number of states $\uparrow\uparrow$

$$p(s', r | s, a)$$



First Visit MC / (Every Visit MC)

- On-policy methods
- Off-policy methods

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

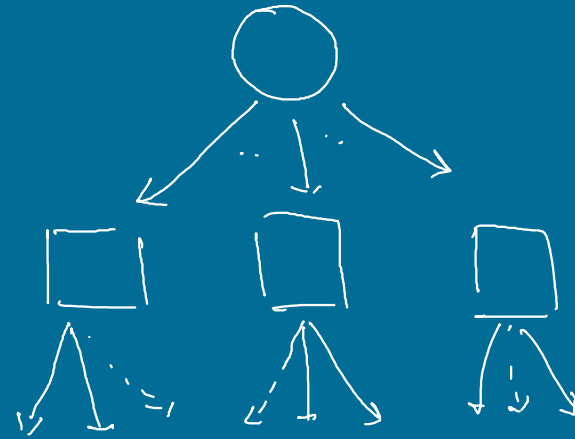
$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

Game Theory

- pick a strategy for player that maximizes his/her utility given the strategies of other players

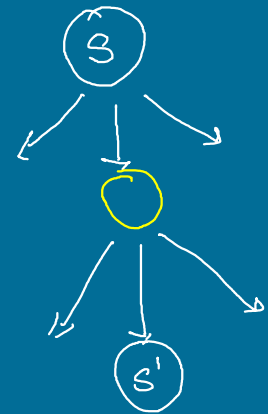
- Study of strategic decision making

Game Tree



UTILITIES

MDP



Probabilistic world outcomes

- John von Neumann, John Nash, Stackelberg

- Rewards

- Returns

Normal Form Game

		1 Actions		
		CRAM	DO-HW-TUT	PLAYGAME
WORLD	EASY	98	100	85
	HARD	97	90	65

$$\text{Utility} \quad U_1(a_{1,k}, \text{world}) \quad a_{1,k} \in A_1$$

$$\text{Exp Utility} = \frac{1}{2} \cdot 98 + \frac{1}{2} \cdot 97$$

$$= p(\text{world} = \text{easy}) \cdot U_1(\text{cram}, \text{easy}) + p(\text{world} = \text{hard}) U_1(\text{cram}, \text{hard})$$

$$U_1(s_1, \text{world}) \quad s_1 - \text{strategy}$$

$$s_1 = (p(a_{1,1}), p(a_{1,2}), \dots, p(a_{1,n}))$$

$$\rightarrow s_1 = (\frac{1}{2} \sim \text{cram}, \frac{1}{2} \sim \text{do-hw-tut}, 0 \sim \text{play game})$$

$$\text{Exp Utility} = p(\text{cram}) \left[p(w = \text{easy}) U_1(\text{cram}, \text{easy}) + p(w = \text{hard}) U_1(\text{cram}, \text{hard}) \right]$$

$$+ p(\text{do-hw-tut}) \left[p(w = e) U_1(\text{do-h-t}, e) + p(w = h) U_1(\text{do-h-t}, h) \right]$$

$$= \frac{1}{2} \left[\frac{1}{2} \cdot 98 + \frac{1}{2} \cdot 97 \right] + \frac{1}{2} \left[\frac{1}{2} \cdot 100 + \frac{1}{2} \cdot 90 \right]$$



(u_1, u_2)

		P2's ACTIONS		
		PLAYER 2		
		ROCK	PAPER	SCISSORS
PLAYER 1	ROCK	0, 0	-1, 1	1, -1
	PAPER	1, -1	0, 0	-1, 1
	SCISSORS	-1, 1	1, -1	0, 0

JOINT UTILITIES

Goal: pick a strategy for player i that maximizes his utility given the strategies of other player.

$$\sum_{i=1}^M u_i(s) = 0$$

↑ strategy

Q. $P_2 \rightarrow \text{Rock}$ $P_1 \rightarrow S_1 = (0, 1, 0)$

$P(R) \quad P(P) \quad P(S)$

• Zero-Sum ~~Profile~~ Game

Q. $P_2 \rightarrow 50\% \text{ Rock} \quad 50\% \text{ Paper}$ $P_1 \rightarrow 50\% \text{ Scissor} \quad 50\% \text{ Paper}$

$E[u_1] = \frac{1}{4}$ $E[u_1] = \frac{1}{2}$

- $u_i(s) = u_i(s_1, s_2, \dots, s_m) = u_i(s_i, s_{-i})$

Payoff/Utility

Alphabet

$$s_{-i} = (s_1, s_2, \dots, s_{i-1}, s_{i+1}, \dots, s_m)$$

Roman

	A	B	C	D	E
i	2,10	4,7	4,6	5,2	3,8
ii	3,8	6,4	5,2	1,3	2,6
iii	5,3	3,1	2,2	4,1	3,0
iv	6,7	9,5	7,5	8,5	5,5

(A, iv) ~ Interesting strategy profile

$$\underline{u_i(s_i, s_{-i})}$$

$$u_{\text{alph}}(A, iv) \geq u_{\text{alph}}(s_{\text{alph}}, iv) \quad \forall s_{\text{alph}} \neq A$$

$$u_{\text{rom}}(iv; A) \geq u_{\text{rom}}(s_{\text{rom}}, A) \quad \forall s_{\text{rom}} \neq iv$$

• Is there always a dominant strategy?

Prisoner's Dilemma

		PRISONER 2	
		Cooperate	Defect
PRISONER 1	Cooperate	-1,-1	-6,0
	Defect	0,-6	-3,-3

$$S^{NE} = (D, D)$$

$$u_1(s_1=D, s_2=D) \geq u_1(\underset{-3}{?}, s_2=D)$$

$$u_2(s_2=D, s_1=D) \geq u_2(\underset{-3}{?}, s_1=D)$$

Social Welfare:

sum of utilities of players

$$SW(C, C) = -1 + (-1) = -2$$

$$SW(D, D) = -3 + (-3) = -6$$

Nash Equilibrium: Nash equilibria are strategy profiles s , where none of the participant benefit from unilaterally changing their decision

$$\{s^{NE} = (s_i^{NE}, s_{-i}^{NE})\} \quad \forall i \quad u_i(s_i^{NE}, s_{-i}^{NE}) \geq u_i(s_i, s_{-i}^{NE}) \quad s_i \neq s_i^{NE}$$

Professors Dilemma

		Student	
		Study	Games
Professor	Effort	1000,1000	0,-10
	Slack	-10,0	0,0

(Effort, Study) (strict)

(Slack, Games) (weak)

Finding Pure Nash Equilibrium

- Find dominating strategy, eliminate all other rows/columns (recurse)
- Remove a strictly dominated strategy (recurse)

Example: 2

	L	C	R
U	10,3	1,5	5,4
M	3,1	2,4	5,2
D	0,10	1,8	7,0

NE
 $S = (M, C)$

Exercise:

	A	B	C	D	E
i	2,4	4,7	4,6	5,2	3,8
ii	3,8	6,4	5,2	1,3	2,6
iii	5,3	3,1	2,2	9,1	3,0
iv	6,7	9,5	5,5	8,5	4,5

Def Weakly strictly dominant strategy equilibrium of a game G ,

$G := (N, A_i, u_i : \prod_i A_i \rightarrow \mathbb{R})$, in strategic form is defined as

the weakly strictly dominant action profile, and is denoted by

$D^W(G)$.

$D^S(G)$

Algo IESD: Iterated Elimination of strictly Weakly Dominated strategies

	L	M	R
U	1,0	1,2	0,1
D	0,3	0,1	2,0

(U, M)

Def A strategic form game is dominance solvable if IESD actions leads to a unique outcome.

	L	R
U	2,1	0,2
D	2,3	4,3

	L	R
U	3,1	2,0
M	4,0	1,1
D	4,4	2,4

(D,R)

Best Response Correspondence

(D,L)

$$B_i : A_{-i} \Rightarrow A_i$$

$$B_i(s_{-i}) = \{ s_i \in A_i : u_i(s_i, s_{-i}) \geq u_i(b_i, s_{-i})$$

	L	M	R
U	1,0	1,2	0,2
D	0,3	1,1	2,0

$$B_1(L) = \{ U \}$$

$$B_1(R) = \{ D \} \quad \forall b_i \in A_i \}$$

$$B_1(M) = \{ U, D \}$$

$$B_2(U) = \{ M, R \}$$

$$B_2(D) = \{ L \}$$

Proposition For any 2-person game in strategic form G ,

$(s_1^*, s_2^*) \in N(G)$ if and only if

$$s_1^* \in B_1(s_2^*) \text{ , } s_2^* \in B_2(s_1^*) \text{ .}$$

	H	T
H	1, -1	-1, 1
T	-1, 1	1, -1

No Nash Equilibrium (pure strategies)

$$\left(\left(\frac{1}{2}, \frac{1}{2} \right), \left(\frac{1}{2}, \frac{1}{2} \right) \right)$$

$$E[u_i] = \frac{\overbrace{P(H)P(H)}^{P(H)P(H)} \cdot u_1(H, H) + \overbrace{P(H)P(T)}^{P(H)P(T)} \cdot u_1(H, T) + \overbrace{P(T)P(H)}^{P(T)P(H)} \cdot u_1(T, H)}{\underbrace{+ P(T)P(T)}_{P(T)P(T)} \cdot u_1(T, T)}$$

$$= 0$$

Definition | A mixed strategy π_i for player i , is a probability distribution over his set of available actions A_i .

$$|A_i| = m \quad \pi_i = (\pi_i^1, \dots, \pi_i^m) \quad \pi_i^k \geq 0, \quad \sum_{k=1}^m \pi_i^k = 1$$

- Let $\Delta(X)$ denote the set of all probability distributions on a set X .

$$\pi_i \in \Delta(A_i)$$

$$u_i(\pi) = \sum_{a \in A} P(a) u_i(a)$$

$$A = \prod_i A_i$$

$$\pi \sim P(a) = \pi_1(a_1) \cdot \pi_2(a_2) \dots \pi_N(a_N)$$

	m	o
m	2, 1	0, 0
o	0, 0	1, 2

$$u_1(p, q) = 2pq + 0 \cancel{p(1-q)} + 0 \cancel{(1-p)q} + (1-p)(1-q)$$

$$\pi = (p, q)$$

$$u_1(p, q) = 1 + 3pq - p - q$$

$$u_2(p, q) = pq + 2(1-p)(1-q) = 2 - 2p - 2q + 3pq$$

$$\# \quad (p, q) := \underline{(1, 1/3)} \quad u_1(1, 1/3) = 2/3 \quad u_2(1, 1/3) = 1/3$$

$$\text{Supp}(\pi_1) = \{m\}$$

$$\text{Supp}(\pi_2) = \{m, o\}$$

Remarks:

- (1) $\pi_i \in B_i(\pi_{-i})$ iff every action in the support of π_i is itself a best response to π_{-i} .
- (2) A mixed strategy profile π^* is MNE iff for each player i , each action in the support of π_i^* is a best response to π_{-i}^* .
- Each action in support of π_i^* yields the same expected payoff (utility) when played against π_{-i}^* , no other action yields a strictly higher payoff.

Proposition | Every finite strategic form game has a mixed strategy equilibrium.

[Kakutani's FPT]

Definition In a strategic form game, player i 's mixed strategy π_i strictly dominates her action a_i^1 if

$$u_i(\pi_i, a_{-i}) > u_i(a_i^1, a_{-i}) \quad \forall a_{-i} \in A_{-i}$$

$p-2$

Example:

$\pi_1 = \left(\overset{T}{0}, \overset{M}{\frac{1}{2}}, \overset{B}{\frac{1}{2}} \right)$ dominates 'T' \rightarrow

$p-1$

$$u_1(\pi_1, L) = \frac{1}{2} \cdot 3 + \frac{1}{2} \cdot 0 = \frac{3}{2} > u_1(T, L)$$

$$u_1(\pi_1, R) = \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 4 = 2 > u_1(T, R)$$

	L	R
T	1,1	1,0
M	3,0	0,3
B	0,1	4,1